# Enhancing Door–Status Detection for Autonomous Mobile Robots during Environment–Specific Operational Use

Michele Antonazzi, Matteo Luperto, Nicola Basilico, N. Alberto Borghese

*Abstract*— Door–status detection, namely recognising the presence of a door and its status (open or closed), can induce a remarkable impact on a mobile robot's navigation performance, especially for dynamic settings where doors can enable or disable passages, changing the topology of the map. In this work, we address the problem of building a door–status detector module for a mobile robot operating in the same environment for a long time, thus observing the same set of doors from different points of view. First, we show how to improve the mainstream approach based on object detection by considering the constrained perception setup typical of a mobile robot. Hence, we devise a method to build a dataset of images taken from a robot's perspective and we exploit it to obtain a door–status detector based on deep learning. We then leverage the typical working conditions of a robot to *qualify* the model for boosting its performance in the working environment via fine–tuning with additional data. Our experimental analysis shows the effectiveness of this method with results obtained both in simulation and in the real–world, that also highlights a trade–off between the costs and benefits of the fine–tuning approach.

## I. INTRODUCTION

Autonomous mobile robots are nowadays increasingly employed for cooperating with humans in a variety of tasks settled in indoor public, private, and industrial workplaces. A challenge posed to these *service robots* is coping with highly dynamic environments characterised by features that can rapidly and frequently change, very often due to the presence of human beings [1]. Consider, as examples, a domestic setup in an apartment or a workspace with several offices. In a time span of hours or days, the topology itself of these environments might frequently change its connectivity, since doors may be left open or closed, hence modifying in time the reachability of free spaces. This phenomenon strongly impacts the capability of robots to efficiently navigate and perform their tasks. At the same time, during their operational time, robots are often exposed to large amounts of data about their surroundings that offer an opportunity to track, model, and predict doors' statuses (and topology variations). The relevance of this problem is well–established in the literature. Different works, such as [2], [3], show how modelling the status of doors across a long time span and predicting the changes in the environment topology improves a robot's task performance. Intuitively, better paths can be planned by taking into account whether a room will be reachable or not upon arriving there.

Central to unlocking such enhanced indoor navigation behaviours is what we call in this work *door–status detection*: the robot's capability to extract, from visual perceptions, the

All authors are with the Department of Computer Science, University of Milan, Milano, Italy `name.surname@unimi.it`
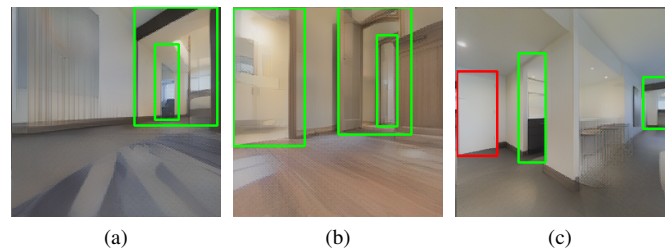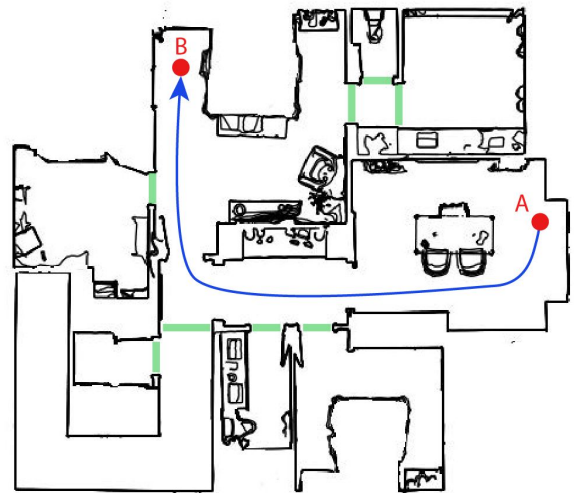
Fig. 1: A robot, navigating from A to B, can observe the status (open or closed) of different doors (highlighted in the top image). In this condition, door–status detection can be a difficult task as, from the robot's point of view, doors can be nested (a–b), doors can be hidden in the wall (c), or instead of a door sometimes there are just passages (a–c). The bounding boxes of open (closed) doors, as identified by our method, are shown in green (red). For the remainder of this work, we follow the same colour schema.

presence and location of a door and, at the same time, to recognise its traversability (open or closed status).

In this work, we propose a method to endow a robot with door–status detection capabilities that can be run during task–related autonomous navigation.

Door–status detection is particularly challenging for mobile robots operating indoor since clear and well–framed views of a door are seldom encountered during navigation. Fig. 1 depicts some typical instances of these challenges. While navigating, the robot can view nested doors (Fig. 1a, 1b), doors that are partially occluded (Fig. 1b, 1c), or closed doors difficult to distinguish from their background (Fig. 1c).

To tackle the above problem, our approach starts with the choice of modelling door–status detection as a variant of object detection (OD) performed with deep neural networks. However, we found that OD deep learning methods, despite their great capabilities, exhibit important shortfalls when cast into the indoor robot navigation setting. Hence, our approach proposes a deployment methodology specific for mobile robots that allows harnessing the potentials of OD based on deep learning while solving what we recognised as the two most important limitations of such techniques in this domain.

First, OD methods are usually trained on large–scale datasets whose images are acquired from a human point of view. As a result, training examples follow a distribution that could be significantly different from the one generating the data perceived by a mobile robot. We show how popular datasets employed to train state–of–the–art deep learning detectors [4], [5], do not properly represent the embodied perception constraints and uncertainty typically characterising a mobile robot [6], thus causing generalisation issues.

Second, deep–learning OD modules are commonly trained with the main objective of obtaining a *general detector*. This model is trained once, stored, and is meant to work in previously unseen environments. These practices are not optimal when considering the typical working conditions of an indoor service robot. After an initial deployment phase, the robot is commonly used in the same environment for a long time, sometimes even for its entire life cycle. In such persistent conditions, the robot eventually observes the same doors multiple times, from different points of view, and under various environmental conditions. Also, different doors may present similar visual features (e.g., multiple doors of the same model). Against this operative background, and from a practical point of view, the ability to generalise in new environments becomes less important, while correctly performing door–status detection in challenging images from the deployment environment becomes paramount.

To address the first limitation, we devise a method for acquiring a large visual dataset from multiple photorealistic simulations taking into account the robot's perception model along realistic navigation paths. This allows us to train a deep *general* door–status detector with examples following a distribution compliant with the robot's perception capabilities. To deal with the second limitation, we exploit the robot's operational conditions to tailor our general detector for a given target environment. We obtain what we call a *qualified detector*, whose performance can substantially improve from the robot's experience enabling door–status detection in challenging instances (see the examples of Fig. 1). Our solution relies on fine–tuning sessions [7]–[9] of the general detector (which shall be considered as a baseline) with new examples from the target environment. These data can be collected and labelled, for example, during the robot installation phases or while the robot carries out its duties. (A setting motivated also by our on–the–field experience with assistive robots [10].)

We evaluate our approach by assessing its performance,

also in the challenging cases exemplified in Fig. 1, with an extensive experimental campaign conducted in simulated settings and in different real–world environments and conditions, as perceived by a mobile robot during its deployment.

## II. RELATED WORKS

Detecting a door's location can be useful for several tasks, as *room segmentation* [11], i.e., to divide the map of the environment into semantically meaningful regions (rooms), to predict the shape of unobserved rooms [12], or to do *place categorisation* [13], [14], which assigns to the rooms identified within the occupancy map a semantic label (e.g., *corridor* or *office*) according to their aspect.

Recent studies [2], [3] show how recognising door statuses can improve the navigation performance of robots in long–term scenarios. The work of [3] models the periodic environmental changes of a dynamic environment in a long–term run, while [2] proposes a navigation system for robots that operate for a long time in indoor environments with traversability changes.

Detecting doors in RGB images has been addressed as an OD task. Classical methods are based on the extraction of handcrafted features [15]–[17]. Deep learning end–to–end methods [18] provide significant improvements thanks to their capability of automatically learning how to characterise an object class, robustly to scale, shift, rotation, and exposure changes. As a significant example, the work of [19] describes a method for door detection with the goal of supporting and improving the autonomous navigation task performed by a mobile robot. A convolutional neural network is trained to detect doors in an indoor environment and its usage is shown to help a mobile robot to traverse passages in a more efficient way. Another approach, proposed in [20], focuses on robustly identifying doors, cabinets, and their respective handles in order to allow grasping by a robot. The authors use a deep architecture based on YOLO [9] to detect the Region Of Interest (ROI) of doors. This allows to obtain the handle's location by focusing only on the area inside the door ROI.

These works are representative examples of methods partially addressing the door–status detection problem in the mobile robotics domain. Indeed they do not explicitly consider the point of view of a mobile robot or do not take advantage of the robot's typical operational conditions. In this work, we devise an approach to overcome such limits.

## III. BUILDING A DOORS DATASET FOR MOBILE ROBOTS

One of the key prerequisites to exploit deep learning to synthesise an effective door–status detector for a mobile robot is the availability of a dataset consistent with its challenging perception model (see Fig. 1). The examples contained in the dataset should follow three main desiderata. Images (i) should represent different environments with different features, thus allowing the model to learn how to generalise; (ii) should contain doors as observed from a point of view similar to the one of a robot navigating in an indoor environment; (iii) should be taken from real environments or with an adequate level of photorealism.

An effective but impractical and time–consuming way to comply with the above requirements would be to deploy a robot on the field and having it exploring different environments while acquiring image samples of doors. The large overheads of such a procedure are well–known and a popular alternative is to rely on simulations [21] or publicly available datasets [4], [22], [23].

Meeting the desiderata (i)-(iii) in simulation is not straightforward since these are seldom guaranteed by available frameworks. For example, simulation tools popular in robotics such as Gazebo [24] or Unreal [25], while providing accurate physics modelling, fail to represent the realism and complexity of the perceptions in the real world. At the same time, public datasets as [4], [22], [23], do not well represent the point of view of a robot in its working conditions [6]. To address these issues, we resorted to Gibson [26], a simulator for embodied agents that focuses on realistic visual perceptions, and to the environments from Matterport3D [27], an RGB–D dataset of 90 real–world scans.

Given a simulated environment, we extract a set of poses that could describe views compatible with a mobile robot by applying a set of principles; the key ones include lying in the reachable free space (feasibility), ensuring a minimum clearance from obstacles, and being along the shortest paths between key connecting locations in the environment's topology. We achieve them with an extraction algorithm working in three phases: grid extraction, navigation graph extraction, and pose sampling.

The grid extraction phase aims at obtaining a 2D occupancy grid map, similar to those commonly used by mobile robots for navigation. We start from the environment's 3D mesh, and we aggregate obstacles from multiple cross–sections of the 3D mesh performed with parallel planes. The result is then manually checked for inaccuracies and artefacts produced during the procedure.

The *navigation graph* (shown in Fig. 2a) is a data structure that we use to represent the topology of the locations on the grid map that correspond to typical waypoints a robot occupies while navigating in the environment. We compute it from a Voronoi tessellation of the grid map by using obstacle cells as basis points [28], extracting graph edges from those locations that maintain maximum clearance from obstacles.

We then perform pose sampling on the navigation graph. The algorithm extracts from the graph a list of positions keeping a minimum distance $D$ between them (this parameter controls the number and the granularity of the samples). A visual example of the algorithm's results is shown in Fig. 2b.

To build the dataset we acquire an image from the points of view of a robot's front–facing camera simulating its perceptions in the virtual environment from the sampled poses. Specifically, in each pose on the grid map, we acquire perceptions at two different height values ($0.1\,\mathrm{m}$ and $0.7\,\mathrm{m}$ – to simulate different embodiments of the robots) and at 8 different orientations (from $0°$ to $315°$ with a step of $45°$). Each acquisition includes the RGB image, the depth information, and the semantic data from Matterport3D. Since the semantic annotation of Matterport3D presents some in-



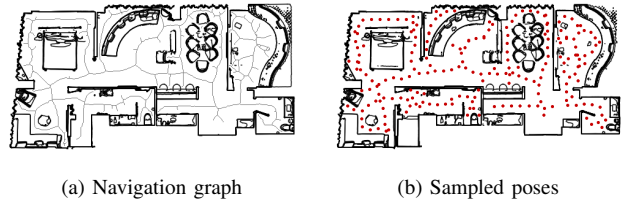(a) Navigation graph     (b) Sampled poses

Fig. 2: Different phases of the pose extraction algorithm.

accuracies, data labelling is manually performed by a human operator who specifies the door bounding boxes and the door status as open or closed. We considered 10 different Matterport3D environments (small apartments or large villas with multiple floors and a heterogeneous furniture style) by setting $D = 1$m. The final dataset we obtained is composed of approx. 5500 examples.

## IV. DOOR–STATUS DETECTION FOR MOBILE ROBOTS

In this section, we first detail how we synthesised a *General Detector* ($GD$, Section IV-A) using a dataset generated with the approach of Section III. Subsequently, leveraging the assumption that the environment $e$ will not change in its core features (location and visual aspect of doors) during the robot's long–term deployment, we introduce our *Qualified Detector* for $e$ ($QD_e$) by applying a procedure based on fine–tuning [7]–[9], [29] of the $GD$ on additional data that, in our envisioned scenario, can be acquired and labelled during the first setup of the robot in $e$.

### A. General Door–Status Detector

As previously introduced, we aim at building and deploying door–status detectors for mobile robots leveraging deep–learning for object detection. In a preliminary experimental phase, we evaluated and compared three popular models suitable for such a task: DETR [29], Faster–RCNN [8], and a YOLO architecture [30]. We decided to adopt DETR since, with respect to the other two methods, it turned out to be easier to deploy in our robotic setting primarily due to two key features. First, DETR does not require setting in advance the number and dimension of anchors (i.e., sets of predefined bounding boxes used to make detections) according to the image resolution and the objects' shape, a task that instead the YOLO architecture requires. Second, both competitors require a final non–maximum suppression step to discard multiple detections of the same object. DETR, instead, matches each bounding box to a different object by construction. Hence, the methodology we describe in this paper and the empirical results evaluating it shall develop around DETR–based detectors. However, we stress the fact that our methods can be applied to any architecture, including those mentioned above, and, eventually, to their improvements.

DETR combines a CNN backbone based on ResNet [7] to produce a compact representation of an image and a transformer [31]. We used the pre–trained version of DETR

on the COCO 2017 [4] dataset and, to adjust for door–status detection, we chose the smallest configuration provided by the authors.

The model requires setting one hyper–parameter, $N$, which determines the fixed number of bounding boxes predicted for each image. As a consequence, to filter out the detected doors, we select the $n \leq N$ bounding boxes whose confidence is not below a threshold $\rho_c$. We tuned $N$ to be higher (but close) to the maximum number of doors in any single image of our dataset.

To train the general detector, we fixed the first two layers of the CNN backbone (as in [29]) with the weights of the pre–trained model. We then re–trained the remaining layers with images from the dataset of Section III. To achieve data augmentation, we generated additional samples by applying a random horizontal flip and resize transformation to a subset of the images (each training sample is selected for this procedure with a probability of 0.5).

### B. Qualification on a Target Environment

Given a new environment $e$ we use a randomly sampled subset of the images collected in it to fine–tune the $GD$, obtaining the qualified detector $QD_e$. To be used in the fine–tuning procedure, these images need to be labelled specifying the bounding boxes and the status for each visible door.

In our envisioned scenario, this data acquisition and labelling tasks can be carried out by a technician during the robot's first installation in $e$ or in a second phase by uploading the data to a remote server. Such a setup phase requires to build the map of the environment (either autonomously or with teleoperation) by observing the entirety of the working environment and is very relevant to many real–world installations of collaborative robots, as we recently experienced with extensive on–the–field testing in the use case of assistive robotics [10]. This manual labelling task is quite time–expensive; yet it is required. In principle, we could use *pseudo–labels* automatically obtained by running the $GD$ over the additional samples for incremental learning. Despite intriguing, we empirically observed that this is particularly challenging due to the fact that pseudo–labels are not enough accurate for this process. Recently, the work of [32] showed how pseudo–labels are particularly noisy and inaccurate: while they can be used to improve performance in tasks where precise labels are less important (like semantic segmentation), they are still too inaccurate to be used in object detection tasks, like the one investigated in this work. We observed how fine–tuning a general detector using pseudo–labels results in a performance degradation of about 20% when compared with the $GD$. These challenges are well-known and the approach we follow in this work is customary. See, for example, the work of [33], where manual annotations have been used to label 3D objects to fine–tune a model employed in long–term localisation. The study of methods to ease the burden of this task will be addressed as a part of our future work.

Finally, note that, while we focus here on a specific $GD$ based on DETR [29], this method to obtain a qualified detector $QD$ is general and can be applied to other deep learning–based models, such as YOLO [30] or Faster–RCNN [8].

## V. EVALUATION IN SIMULATION

### A. Experimental Setting

We evaluate our method using simulated data $\mathcal{D}$ obtained, as described in Section III, from 10 different Matterport3D [27] environments. We test the performance of our detectors on each environment $e$ independently. First, we train the general detector $GD_{-e}$ using the dataset $\mathcal{D}_{-e}$, where $\mathcal{D} = \{\mathcal{D}_{-e}, \mathcal{D}_e\}$, $\mathcal{D}_e$ contains all the instances acquired from poses sampled in environment $e$, and $\mathcal{D}_{-e} = \mathcal{D} \setminus \mathcal{D}_e$. This general detector will be used as a baseline in most of the evaluations we present. Then, we randomly partition the first subset as $\mathcal{D}_e = \{\mathcal{D}_{e,1}, \mathcal{D}_{e,2}, \mathcal{D}_{e,3}, \mathcal{D}_{e,4}\}$, where each $\mathcal{D}_{e,i}$ contains the 25% of the examples from $e$, randomly selected.

While $\mathcal{D}_{e,4}$ is reserved for testing, the remaining subsets are used to perform a series of fine–tuning rounds to obtain the corresponding qualified door–status detectors. Specifically, we fine–tune $GD_{-e}$ using these three additional data subsets: $\{\mathcal{D}_{e,1}\}$, $\{\mathcal{D}_{e,1}, \mathcal{D}_{e,2}\}$, and $\{\mathcal{D}_{e,1}, \mathcal{D}_{e,2}, \mathcal{D}_{e,3}\}$. We denote the obtained qualified detectors as $QD_e^{25}$, $QD_e^{50}$, and $QD_e^{75}$, respectively. The superscript denotes the percentage of data instances from $e$ that are required to fine–tune the general door–status detector. Such a percentage can be interpreted as an indicator of the cost to acquire and label the examples. To give a rough idea, labelling the 25% of the dataset (approximately 150 images) took a single human operator an effort of about 1 hour.

We empirically set the parameters of the door–status detector as $N = 10$ (number of bounding boxes) and $\rho_c = 0.75$ (confidence threshold). We conducted an extensive preliminary experimental campaign spanning different batch sizes ($\{1, 2, 4, 16, 32\}$) and the number of epochs ($\{20, 40, 60\}$) selecting 1 and 60 for the general detector and 1 and 40 for the qualified ones, respectively.

We measure performance with the average precision score (AP) used in the Pascal VOC challenge [5] by adjusting for a finer interpolation of the precision/recall curve to get a more conservative (in the pessimistic sense) evaluation. The AP is a popular evaluation metric widely adopted for object detection tasks, it represents the shape of the precision/recall curve as the mean precision over evenly distributed levels of recall. To accept a true positive, the bounding box computed by the network must exhibit an Intersection Over Union area ($IOU$) with one true bounding box above a threshold $\rho_a$, that we empirically set to 50%.

The source code of our simulation framework (Section III), the door–status detectors (Section IV), and the collected datasets are maintained in a freely accessible repository[1].

### B. Results

Table I reports the mean AP scores (averaged over the 10 environments) reached by the 4 detectors divided by label

---

[1] https://github.com/aislabunimi/
door-detection-long-term

(closed door, and open door), the average increments (with respect to the detector immediately above in table) obtained with fine–tuning, and the standard deviation ($\sigma$). We also report the AP scores for every environment in Fig. 3. These results show the trade-off between performance increase (via fine–tuning) and costs due to data collection and labelling.

| Exp. | Label | AP | $\sigma$ | Increment | $\sigma$ |
|------|-------|----|---|----------|----|
| $GD_{-e}$ | Closed | 34 | 12 | – | – |
| | Open | 48 | 12 | – | – |
| $QD_e^{25}$ | Closed | 55 | 15 | 70% | 58 |
| | Open | 60 | 10 | 30% | 34 |
| $QD_e^{50}$ | Closed | 64 | 12 | 21% | 21 |
| | Open | 68 | 10 | 14% | 11 |
| $QD_e^{75}$ | Closed | 72 | 10 | 14% | 9 |
| | Open | 72 | 9 | 7% | 5 |

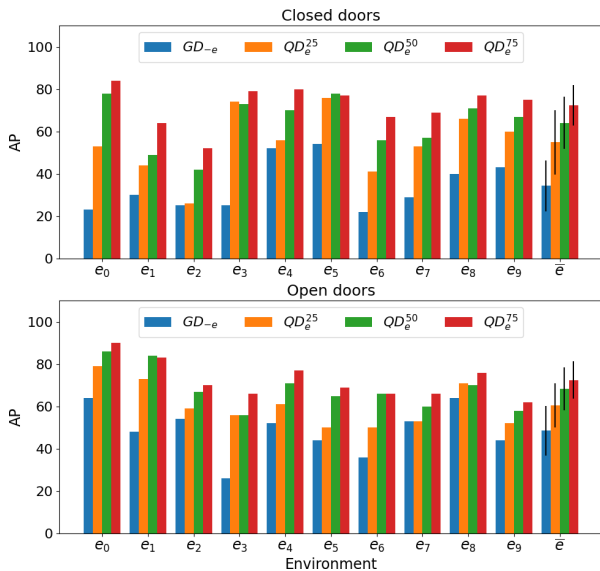TABLE I: Average AP in Matterport3D environments.



Fig. 3: AP scores in Matterport3D environments.

Results from Table I and Fig. 3 show that the general detector $GD_{-e}$, thanks to our dataset's consistency with the robot's perception model (see Section III), is able to correctly detect doors statuses in those cases where they are clearly visible, as shown in Fig. 4. However, while such a performance allows its use on a robot, there is significant room for improvement in detecting more challenging examples.
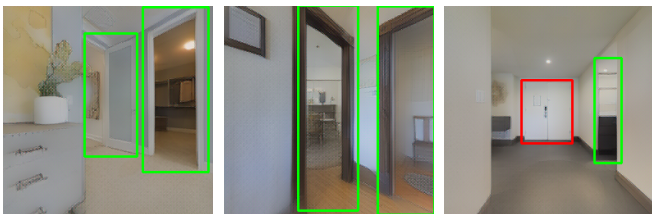


Fig. 4: Door–statuses found by the general detector $GD_{-e}$.

More interestingly, qualified detectors achieve a steep increase in performance. Unsurprisingly, the performance improves with more data (and data preparation costs) from $QD_e^{25}$ to $QD_e^{75}$. However, it can be seen how $QD_e^{25}$, despite requiring a relatively affordable effort in manual labelling, obtains the highest performance increase. From a practical perspective, this shows how the availability of *a few* labelled examples from the robot target environment could be a good compromise between performance and costs to develop an environment–specific door–status detector. This suggests that the number of examples that have to be collected and labelled on the field can be limited, thus promoting the applicability of our proposed framework. An example of this is shown in Fig. 5, where it can be seen how a $QD_e^{25}$ (bottom row) fixes the mistakes of its corresponding general detector $GD_{-e}$ (top row) in challenging images with nested or partially observed doors.



Fig. 5: Door–statuses as identified by $GD_{-e}$ (top row) compared to $QD_e^{25}$ (bottom row) in Matterport3D environments.

## VI. EVALUATION IN THE REAL WORLD

### A. Experimental Settings

In this section, we evaluate the performance of our method with a real robot by using images collected by a Giraff–X platform [10] (Fig. 6a) during a teleoperated exploration of 3 single–floor indoor environments with multiple rooms from two buildings in our campus. Images were extracted from the robot's perceptions during navigation at 1 fps. As it commonly happens with real–world robot data, images are acquired in noisy environmental conditions with low–quality cameras, thus making the detection task even more difficult. In our setting, we used 320x240 RGB images acquired with an Orbbec Astra RGB–D camera.

First, we consider two general detectors. One is trained with simulated data $\mathcal{D}$, as in the previous section but with all the 10 environments. Another one is trained with real–world images from the publicly available *DeepDoors2* dataset (DD2) [23], which features 3000 images of doors that we re-labelled to include the ground truth for challenging examples not originally provided. Comparing these two general detectors, we aim at assessing the advantages of training with a dataset following the principles we proposed in Section III

instead of relying on mainstream datasets for classical object detection.

Subsequently, following the same steps of Section V-A, we qualify both $GD$s by using the 25, 50, and 75% of data collected in the three real environments.

## B. Evaluation metrics

Analogously to what is done in simulation, we report the AP scores, but we argue that the real–world evaluation of our method can be conducted also with additional metrics, which are more representative of the actual application domain where door–status detection is meant to be cast.

The AP (as well as other metrics used in computer vision) presents some limitations when used in our context. Such a metric considers as false positives multiple bounding boxes assigned to the same door, as in Fig. 6b. However, a robot can easily disambiguate this by leveraging additional data such as its estimated pose and the map of the environment. On one side, although an erroneous localisation of the bounding box of a door (Fig. 6c) penalises the AP, it might have little effect in practice. On the other side, the AP is very marginally affected by a wrong detection of the door status if the bounding box is sufficiently accurate due to the fact that different labels are treated as two independent object classes, as the case in Fig. 6b. Conversely, the error of misleading a closed passage for an open one (and *vice versa*) can significantly impact the robot's performance when the robot translates such information into actions.
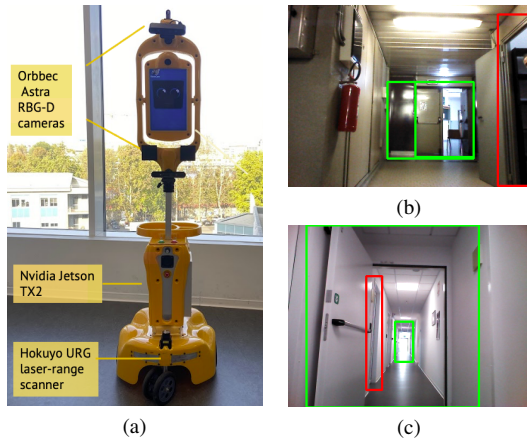


Fig. 6: Examples (b–c) of different types of errors made by a door–status detector mounted on our Giraff–X robot (a). While the AP considers all these errors in a similar way, our proposed metric considers them differently, according to their potential impact on robot performance.

To address this shortcoming and better capture performance in our robotic setting, we introduce three additional metrics. Consider a door $i$ in a given image. If multiple bounding boxes are matched to $i$, where matching means $IOU \geq \rho_a$, the one with the maximum above–threshold ($\rho_c$) confidence is selected. If the status of the door is correctly identified, we consider it as a true positive ($TP$). Otherwise, we classify it as a False Positive ($FP$). All the remaining

bounding boxes matched to $i$ are, as per our previous considerations, ruled out from the evaluation. Finally, when a bounding box does not meet the $IOU$ condition with any door in the image, we count it as a Background False Detection ($BFD$). A False Positive and a Background False Detection are errors that can play very different roles inside a robotic use case. While the first is likely to affect the robot's decisions, the second one might increase the uncertainty in the robot world–model. We scale the above metrics using the true number of doors in the testing set, denoted as $GT$, thus obtaining $TP_\% = TP/GT$, $FP_\% = FP/GT$, and $BFD_\% = BFD/GT$.

## C. Results

Table II compares the average AP of the $GD$ trained with DD2 [23] with that trained with our dataset $\mathcal{D}$. Intuitively, a model trained with real–world data (such as those featured in DD2) should have higher performance when used with real–world images, if compared with a model trained with simulated data (as $\mathcal{D}$). However, Table II shows how the $GD$ and $QD$s trained with $\mathcal{D}$ have higher performance than those trained with DD2. This is because training images of $\mathcal{D}$, collected from the simulated point of view of a robot, better represent the actual distribution of robot perceptions, allowing us to fill, to some extent, the sim–to–real gap. Moreover, Table II shows that the fine–tuning operation to qualify general detectors to the target environment works remarkably well also when used in real–world conditions. For these reasons, from now on, we present results referring to the general detector trained with $\mathcal{D}$.

| Exp. | Label | DeepDoors2 (DD2) | | | | Simulation dataset ($\mathcal{D}$) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | AP | $\sigma$ | Inc. | $\sigma$ | AP | $\sigma$ | Inc. | $\sigma$ |
| $GD$ | Closed | 5 | 3 | – | – | 13 | 10 | – | – |
| | Open | 18 | 5 | – | – | 31 | 11 | – | – |
| $QD_e^{25}$ | Closed | 33 | 9 | 631% | 240 | 53 | 9 | 508% | 424 |
| | Open | 43 | 14 | 134% | 20 | 55 | 14 | 83% | 19 |
| $QD_e^{50}$ | Closed | 52 | 7 | 66% | 51 | 65 | 8 | 24% | 15 |
| | Open | 51 | 14 | 18% | 7 | 70 | 7 | 29% | 22 |
| $QD_e^{75}$ | Closed | 55 | 8 | 5% | 6 | 72 | 8 | 10% | 5 |
| | Open | 65 | 8 | 32% | 18 | 78 | 8 | 13% | 1 |

TABLE II: Average AP in real–world environments when DD2 dataset and our one ($\mathcal{D}$) are used to train the $GD$.

The performance of $QD$s is similar to that obtained in the far less challenging dataset of Table I. Most importantly, the performances of $QD_e^{25}$ corroborate our findings from Section V-A, confirming how few additional examples can induce a significant increase in performance. Beyond the improvements observed in the average scores, $QD_e^{25}$ managed to provide correct door–status detection in very challenging cases. We report in Fig. 7 some representative examples. As it can be seen, our detectors correctly recognised cases with nested doors, partially visible frames, and narrow side views. Another relevant example is in the second image (top row) of Fig. 7, where the qualified model succeeds in detecting a white closed door in the background while, at the same time, not making a false detection of the white wardrobe doors on the right.

Fig. 7: Challenging door–statuses detected by $QD_e^{25}$ in the real environments $e_1$, $e_2$, and $e_3$ (ordered by columns).



Fig. 8: $TP_\%$, $FP_\%$, and $BFD_\%$ obtained by $QD_{e_3}^{25}$ for increasing confidence thresholds.

| Env. | Exp. | GT | TP (TP%) | FP (FP%) | BFD (BFD%) |
|------|------|----|----------|----------|------------|
| $e_1$ | $GD$ | 235 | 71 (30%) | 18 (7%) | 51 (21%) |
| | $QD_e^{25}$ | 235 | 145 (61%) | 10 (4%) | 64 (27%) |
| | $QD_e^{50}$ | 235 | 179 (76%) | 4 (1%) | 44 (18%) |
| | $QD_e^{75}$ | 235 | 190 (80%) | 4 (1%) | 36 (15%) |
| $e_2$ | $GD$ | 269 | 96 (35%) | 17 (6%) | 56 (20%) |
| | $QD_e^{25}$ | 269 | 192 (71%) | 11 (4%) | 87 (32%) |
| | $QD_e^{50}$ | 269 | 206 (76%) | 6 (2%) | 66 (24%) |
| | $QD_e^{75}$ | 269 | 228 (84%) | 7 (2%) | 60 (22%) |
| $e_3$ | $GD$ | 327 | 62 (18%) | 19 (5%) | 108 (33%) |
| | $QD_e^{25}$ | 327 | 183 (55%) | 22 (6%) | 190 (58%) |
| | $QD_e^{50}$ | 327 | 230 (70%) | 13 (3%) | 103 (31%) |
| | $QD_e^{75}$ | 327 | 248 (75%) | 8 (2%) | 75 (22%) |

TABLE III: Extended results in the real–world environments.

| Exp. | Label | AP | $\sigma$ | Increment | $\sigma$ |
|------|-------|----|----|-----------|----|
| $GD$ | Closed | 14 | 18 | – | – |
| | Open | 31 | 8 | – | – |
| $QD_e^{25}$ | Closed | 38 | 8 | 781% | 1016 |
| | Open | 45 | 11 | 46% | 3 |
| $QD_e^{50}$ | Closed | 48 | 8 | 26% | 8 |
| | Open | 53 | 17 | 17% | 8 |
| $QD_e^{75}$ | Closed | 54 | 8 | 14% | 1 |
| | Open | 56 | 16 | 6% | 5 |

TABLE IV: Average AP results in $e_1$ and $e_2$ tested in different light conditions with respect to those used to qualify the detector (day/night time).

| Env. | Exp. | GT | TP (TP%) | FP (FP%) | BFD (BFD%) |
|------|------|----|----------|----------|------------|
| $e_1$ | $GD$ | 1079 | 334 (30%) | 56 (5%) | 150 (13%) |
| | $QD_e^{25}$ | 1079 | 532 (49%) | 62 (5%) | 306 (28%) |
| | $QD_e^{50}$ | 1079 | 572 (53%) | 65 (6%) | 299 (27%) |
| | $QD_e^{75}$ | 1079 | 634 (58%) | 56 (5%) | 248 (22%) |
| $e_2$ | $GD$ | 1051 | 335 (31%) | 68 (6%) | 276 (26%) |
| | $QD_e^{25}$ | 1051 | 584 (55%) | 55 (5%) | 357 (33%) |
| | $QD_e^{50}$ | 1051 | 690 (65%) | 40 (3%) | 217 (20%) |
| | $QD_e^{75}$ | 1051 | 700 (66%) | 48 (4%) | 236 (22%) |

TABLE V: Extended results in $e_1$ and $e_2$ tested in different light conditions with respect to those used to qualify the detector (day/night time).

Table III reports the detailed results, for all three environments, of the metrics we defined in Section VI-B. The results show that $GD$, although it has a low number of wrong predictions ($FP$ and $BFD$), is capable of detecting only a few of the $GT$ doors in the images ($TP$). On the contrary, $QD$s dramatically improve performance, with $QD_e^{25}$ showing a $TP_\%$ of 62% on average.

Among the three environments considered, we argue that $e_3$ is the more challenging, as it can be seen by the higher number of $BFD$. To cope with this, there are two possible directions. First, increasing the number of manually labelled examples reduces $BFD$ (as can be seen already with $QD_e^{50}$). Alternatively, adopting a more conservative selection rule by increasing the confidence threshold $\rho_c$, at the cost of slightly reducing the number of $TP$. In Fig. 8, we show how $TP_\%$, $FP_\%$, and $BFD_\%$ for $QD_{e_3}^{25}$ change when varying $\rho_c$ in $e_3$. Such an instance confirms how $\rho_c = 0.75$ is an acceptable trade–off among $TP_\%$ (high) and $BFD_\%$ (low) for such a detector.

In long–term runs, the illumination conditions of an environment might change from those of the initial setup, and this may affect the performance of the door–status detector. To test the robustness of our approach to this event, we acquire (following the same procedure of Section VI-A) data from environments $e_1$ and $e_2$ during nighttime, when only artificial light is present and some rooms are dark. Then, we use these images to test the $GD$ and the $QD$s fine-tuned with data acquired during the initial setup time, with daylight.

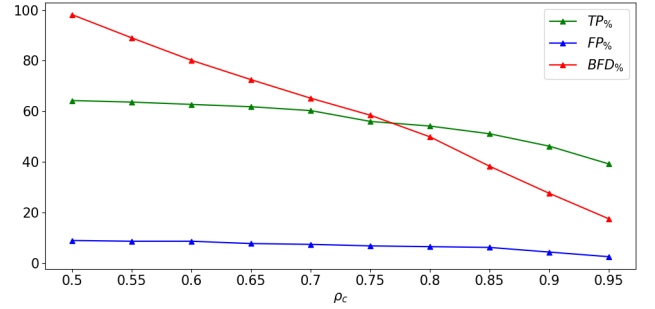The average AP obtained with different lighting conditions is reported in Table IV while the results of our extended metric (presented in Section VI-B) are shown in Table V. Comparing them with Tables II and III respectively, we can see how the performances of the $GD$ are robust to illumination changes, as they are similar to those obtained during daytime. More interestingly, it can be seen how the improvement of $QD$s from the fine–tune is maintained also with different light conditions, with a slight performance decrease if compared to the results of Tables II and III. This is a direct consequence of the fine–tune, which produces $QD$s that slightly overfit the illumination conditions seen during training. Despite this, our method ensures a performance improvement to the $GD$ when used in long–term scenarios with illumination changes, enabling the $QD$s to still solve challenging examples, as shown in Fig. 9. Once again, $QD_e^{25}$, albeit using a few examples for fine–tuning, ensures the best performance improvement also under variable light conditions. See the video attachment, also linked in the repository, for additional examples of our method.
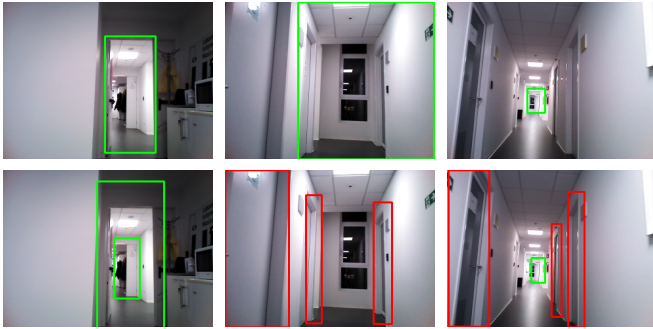
Fig. 9: Challenging nighttime examples classified by $GD$ (top row) and $QD_e^{25}$ (bottom row). $QD_e^{25}$ is fine–tuned with examples obtained with daylight.

## VII. CONCLUSIONS

In this work, we presented a door–status detection method for mobile robots. Our method, based on a deep learning architecture, allows robots to recognise open or closed doors in challenging situations. To train our model, we built a dataset of labelled images from photorealistic simulations taking into account the point of view of a mobile robot. We then fine–tuned a general model into a qualified one to increase performance in the robot's working environment.

Future work will investigate how to quantify and reduce the effort needed for labelling examples to qualify a general detector. Furthermore, we will investigate online fine–tuning methods towards the goal to have a robot that can learn with experience to better distinguish features in its environment.

## REFERENCES

[1] L. Kunze, N. Hawes, T. Duckett, M. Hanheide, and T. Krajnik, "Artificial Intelligence for Long-Term Robot Autonomy: A Survey," *IEEE RA-L*, vol. 3, no. 4, pp. 4023–4030, 2018.

[2] L. Nardi and C. Stachniss, "Long-term robot navigation in indoor environments estimating patterns in traversability changes," in *Proc. ICRA*, 2020, pp. 300–306.

[3] T. Krajník, J. P. Fentanes, J. M. Santos, and T. Duckett, "Fremen: Frequency map enhancement for long-term mobile robot autonomy in changing environments," *IEEE Trans. Rob.*, vol. 33, no. 4, pp. 964–977, 2017.

[4] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common objects in context," in *Proc. ECCV*, 2014, pp. 740–755.

[5] M. Everingham, L. V. Gool, C. K. I. Williams, J. M. Winn, and A. Zisserman, "The pascal visual object classes (VOC) challenge," *Int. J. Comput. Vision*, vol. 88, pp. 303–338, 2009.

[6] N. Sünderhauf, O. Brock, W. Scheirer, R. Hadsell, D. Fox, J. Leitner, B. Upcroft, P. Abbeel, W. Burgard, M. Milford, and P. Corke, "The limits and potentials of deep learning for robotics," *Int. J. Rob. Res.*, vol. 37, no. 4-5, pp. 405–420, 2018.

[7] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. CVPR*, 2016.

[8] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 6, pp. 1137–1149, 2017.

[9] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proc. CVPR*, 2016.

[10] M. Luperto, M. Romeo, J. Monroy, J. Renoux, A. Vuono, F.-A. Moreno, J. Gonzalez-Jimenez, N. Basilico, and N. A. Borghese, "User feedback and remote supervision for assisted living with mobile robots: A field study in long-term autonomy," *Robot Auton Syst*, vol. 155, p. 104170, 2022.

[11] R. Bormann, F. Jordan, W. Li, J. Hampp, and M. Hägele, "Room segmentation: Survey, implementation, and analysis," in *Proc. ICRA*, 2016, pp. 1019–1026.

[12] M. Luperto, F. Amadelli, M. Di Berardino, and F. Amigoni, "Mapping beyond what you can see: Predicting the layout of rooms behind closed doors," *Robot Auton Syst*, vol. 159, p. 104282, 2023.

[13] P. Espinace, T. Kollar, A. Soto, and N. Roy, "Indoor scene recognition through object detection," in *Proc. ICRA*, 2010, pp. 1406–1413.

[14] N. Sünderhauf, F. Dayoub, S. McMahon, B. Talbot, R. Schulz, P. Corke, G. Wyeth, B. Upcroft, and M. Milford, "Place categorization and semantic mapping on a mobile robot," in *Proc. ICRA*, 2016, pp. 5729–5736.

[15] N. Kwak, H. Arisumi, and K. Yokoi, "Visual recognition of a door and its knob for a humanoid robot," in *Proc. ICRA*, 2011, pp. 2079–2084.

[16] I. Monasterio, E. Lazkano, I. Rañó, and B. Sierra, "Learning to traverse doors using visual information," *Math. Comput. Simul.*, vol. 60, no. 3, pp. 347–356, 2002.

[17] X. Yang and Y. Tian, "Robust door detection in unfamiliar environments by combining edge and corner features," in *Proc. CVPR*, 2010, pp. 57–64.

[18] A. Voulodimos, N. Doulamis, A. Doulamis, and E. Protopapadakis, "Deep learning for computer vision: A brief review," *Comput. Intell. Neurosci.*, vol. 2018, 2018.

[19] W. Chen, T. Qu, Y. Zhou, K. Weng, G. Wang, and G. Fu, "Door recognition and deep learning algorithm for visual based robot navigation," in *Proc. ROBIO*, 2014, pp. 1793–1798.

[20] A. Llopart, O. Ravn, and N. A. Andersen, "Door and cabinet recognition using convolutional neural nets and real-time method fusion for handle detection and grasping," in *Proc. ICCAR*, 2017, pp. 144–149.

[21] J. Collins, S. Chand, A. Vanderkop, and D. Howard, "A review of physics simulators for robotic applications," *IEEE Access*, vol. 9, pp. 51 416–51 431, 2021.

[22] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proc. CVPR*, 2009, pp. 248–255.

[23] J. Ramôa, V. Lopes, L. Alexandre, and S. Mogo, "Real-time 2d–3d door detection and state classification on a low-power device," *SN Appl. Sci.*, vol. 3, 2021.

[24] N. Koenig and A. Howard, "Design and use paradigms for gazebo, an open-source multi-robot simulator," in *Proc. IROS*, 2004, pp. 2149–2154.

[25] S. Carpin, M. Lewis, J. Wang, S. Balakirsky, and C. Scrapper, "Usarsim: a robot simulator for research and education," in *Proc. ICRA*, 2007, pp. 1400–1405.

[26] F. Xia, A. R. Zamir, Z. He, A. Sax, J. Malik, and S. Savarese, "Gibson env: Real-world perception for embodied agents," in *Proc. CVPR*, 2018, pp. 9068–9079.

[27] A. X. Chang, A. Dai, T. A. Funkhouser, M. Halber, M. Nießner, M. Savva, S. Song, A. Zeng, and Y. Zhang, "Matterport3d: Learning from RGB-D data in indoor environments," *arXiv:1709.06158*, 2017.

[28] S. Thrun and A. Bücken, "Integrating grid-based and topological maps for mobile robot navigation," in *Proc. AAAI*, 1996, pp. 944–951.

[29] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *Proc. ECCV*, 2020.

[30] A. Farhadi and J. Redmon, "YoloV3: An incremental improvement," in *Proc. CVPR*, 2018, pp. 1804–2767.

[31] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. NIPS*, 2017, pp. 5998–6008.

[32] J. Frey, H. Blum, F. Milano, R. Siegwart, and C. Cadena, "Continual adaptation of semantic segmentation using complementary 2d-3d data representations," *IEEE-RAL*, vol. 7, no. 4, pp. 11 665–11 672, 2022.

[33] N. Zimmerman, T. Guadagnino, X. Chen, J. Behley, and C. Stachniss, "Long-term localization using semantic cues in floor plan maps," *IEEE Robotics and Automation Letters*, vol. 8, no. 1, pp. 176–183, 2023.