

# Spot the Difference: A Novel Task for Embodied Agents in Changing Environments

Federico Landi, Roberto Bigazzi, Marcella Cornia, Silvia Cascianelli, Lorenzo Baraldi, Rita Cucchiara  
University of Modena and Reggio Emilia  
Email: {name.surname}@unimore.it

**Abstract**—Embodied AI is a recent research area that aims at creating intelligent agents that can move and operate inside an environment. Existing approaches in this field demand the agents to act in completely new and unexplored scenes. However, this setting is far from realistic use cases that instead require executing multiple tasks in the same environment. Even if the environment changes over time, the agent could still count on its global knowledge about the scene while trying to adapt its internal representation to the current state of the environment. To make a step towards this setting, we propose *Spot the Difference*: a novel task for Embodied AI where the agent has access to an outdated map of the environment and needs to recover the correct layout in a fixed time budget. To this end, we collect a new dataset of occupancy maps starting from existing datasets of 3D spaces and generating a number of possible layouts for a single environment. This dataset can be employed in the popular Habitat simulator and is fully compliant with existing methods that employ reconstructed occupancy maps during navigation. Furthermore, we propose an exploration policy that can take advantage of previous knowledge of the environment and identify changes in the scene faster and more effectively than existing agents. Experimental results show that the proposed architecture outperforms existing state-of-the-art models for exploration on this new setting.

## I. INTRODUCTION

Imagine you have just bought a personal robot, and you ask it to bring you a cup of tea. It will start roaming around the house while looking for the cup. It probably will not come back until some minutes, as it is new to the environment. After the robot knows your house, instead, you expect it to perform navigation tasks much faster, exploiting its previous knowledge of the environment while adapting to possible changes of objects, people, and furniture positioning. Embodied AI has recently gained a lot of attention from the research community, with amazing results in challenging tasks such as visual exploration [1], [2], [3] and navigation [4], [5], [6], [7]. However, in the current setting, the environment is completely unknown to the agent as a new episode begins. We believe that this choice is not supported by real-world experience, where information about the environment can be stored and reused for future tasks. As agents are likely to stay in the same place for long periods, such information may be outdated and inconsistent with the actual layout of the environment. Therefore, the agent also needs to discover those differences during navigation. In this paper, we introduce a new task for Embodied AI, which we name *Spot the Difference*. In the proposed setting, the agent must identify all the differences between an outdated map of the environment and its current state – a challenge

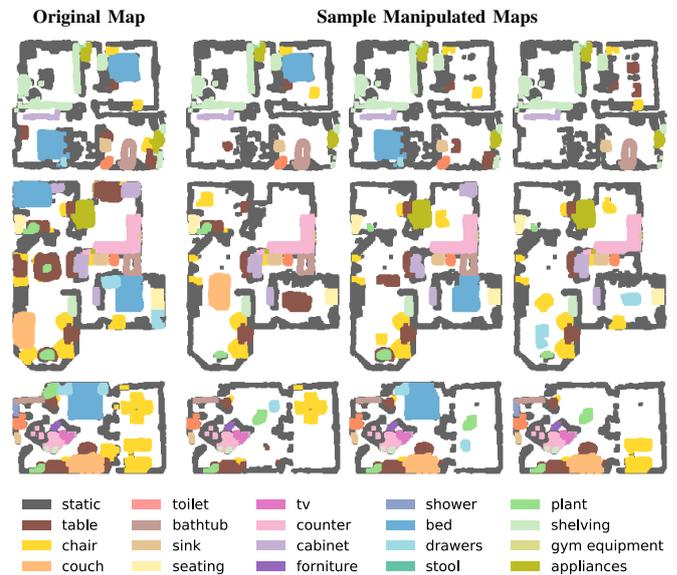


Fig. 1. Generation of alternative states of an environment: original and sample manipulated semantic maps.

that combines visual exploration using monocular images and embodied navigation with spatial reasoning. To succeed in this task, the agent needs to develop efficient exploration policies to focus on likely changed areas while exploiting priors about objects of the environment. We believe that this task could be useful to train agents that will need to deal with changing environments.

Recent work on Embodied AI has tackled the training of embodied agents capable of navigating and locating objects [8], [5], [6], [9]. One of the key factors for success in the field consists in building map representations in which knowledge about the environment can be stored while the agent proceeds [10], [4]. However, the dominant training and evaluation protocol involves an agent initialized from scratch that sees the environment for the first time [11]. Another line of work [12], [13], [1], [14], [15], [3], instead, introduces a mapping phase of the environment to increase the performance on both exploration and down-stream tasks. Unfortunately, if the environment changes over time, the agent needs to rebuild a full representation from scratch and cannot count on an efficient policy to update its internal representation of the environment. In this work, we simulate the natural evolution of an environment and design a specific policy to navigate in

changing environments seamlessly.

Due to the high cost of 3D acquisitions from the real world, the existing datasets of photorealistic 3D spaces [16], [17] do not contain different layouts for the same environment. In this paper, we create a reproducible set-up to generate alternative layouts for an environment. We semi-automatically build a dataset of 2D semantics occupancy maps in which the objects can be removed and rearranged while the area and the position of architectural elements do not change (Fig. 1). In the proposed setting, the agent is deployed in an interactive 3D environment and provided with a map from our produced dataset, which represents the information retained while performing tasks in a past state of the environment.

To train agents that can deal with changing environments efficiently, we develop a novel reward function and an approach for navigation aiming at finding relevant differences between the previous layout of the environment and the current one. Our method is based on the recent Active Neural SLAM paradigm proposed in [10] and [4]. Differently from previous proposals, though, it can read and update the given map to identify relevant differences in the environment in the form of their projections on the map. Our dataset and architecture can be employed with the Habitat simulator [18], a popular research platform for Embodied AI that renders photorealistic scenes and that enables seamless sim-to-real deployment of navigation agents [19], [20]. Experimental results show that our approach performs better than existing state-of-the-art architectures for exploration in our newly-proposed task. We also compare with different baselines and evaluate our agent in terms of percentage of area seen, percentage of discovered differences, and metric curves at varying exploration time budgets. The new dataset, together with our code and pretrained models, is available at this link.

## II. RELATED WORK

Current research directions on Embodied AI for navigation agents can be categorized according to the quantity of knowledge about the environment provided to the agent prior to performing the task [11]. The first direction focuses on the scenario in which the agent is deployed in a completely new environment for which it has no prior knowledge [21], [10], [14]. Running exploration in parallel with a target-driven navigation task resulted in an effective approach to solve the latter (*e.g.*, object-goal navigation [5] and point-goal navigation [4]). Other directions consider the case in which the agent can exploit pre-acquired information about the environment [22], [23] when performing a navigation task. Such pre-acquired information can be either partial [24], [25], [26] or complete [12], [8], [1]. A major limitation of such approaches is that the obtained map representation is assumed to conform perfectly with the environment where the downstream task will be performed.

In this work, we explore a fourth direction, in which the pre-acquired map provided to the agent is incomplete or incorrect due to changes occurred in the environment over time. Common strategies to deal with changing environments

TABLE I  
NUMBER OF MANIPULATED MAPS GENERATED PER DATASET SPLIT.

Dataset Split	Semantic Classes	Scans	Generated SOMs	Episodes
MP3D Train	42	58	2070	$\approx 4.5\text{M}$
MP3D Validation	42	9	160	320
MP3D Test	42	14	260	610
Gibson Validation	20	5	130	450

entail disregarding dynamic objects as landmarks when performing SLAM [27], [28] and applying local policies to avoid them when navigating [29]. An alternative strategy is learning to predict geometric changes based on experience, as done in [30], where the environment is represented as a traversability graph. The main limitation of this strategy is its computational intractability when considering dense metric maps of wide areas, as in our case.

## III. PROPOSED SETTING

In the first part of this section, we introduce a new task for embodied agents, named *Spot the Difference*. We then describe the newly-proposed dataset that we create to enable this setting. Finally, we propose an architecture for embodied agents to tackle the defined task.

### A. *Spot the Difference*: Task Definition

At the beginning of an episode, the agent is spawned in a 3D environment and is given a pre-built occupancy map  $M$ , representing its spatial knowledge of the environment, *i.e.* a previous state of the environment that is now obsolete:

$$M = (m_{ij}) \in [0, 1], \quad 0 \leq i, j < W, \quad (1)$$

where  $m_{ij}$  represents the probability of finding an obstacle at coordinates  $(i, j)$ . The task entails exploring the current environment to recognize all the differences with respect to the state in which  $M$  was computed, in the form of free and occupied space. To accomplish the task, the agent should build a correct occupancy map of the current environment starting from  $M$ , recognizing and focusing on parts that are likely to change (*e.g.*, the middle of wide rooms rather than tight corridors).

For every episode of *Spot the Difference*, the agent is given a time budget of  $T$  time-steps. At time  $t = 0$ , the agent holds the starting map representation  $M$ . At each time-step  $t$ , the map is updated depending on the current observation to obtain  $M_t$ . Whenever the agent discovers a new object or a new portion of free space, the internal representation of the map changes accordingly. The goal is to gather as much information as possible about changes in the environment by the end of the episode. To measure the agent performance, we compare the final map  $M_T$  produced by the agent with the ground-truth occupancy map  $M^*$ . In this sense, the paradigm we adopt is the one of knowledge reuse starting from partial knowledge.

### B. Dataset Creation

**Semantic Occupancy Map.** Given a 3D environment, we place the agent in a free navigable location with heading

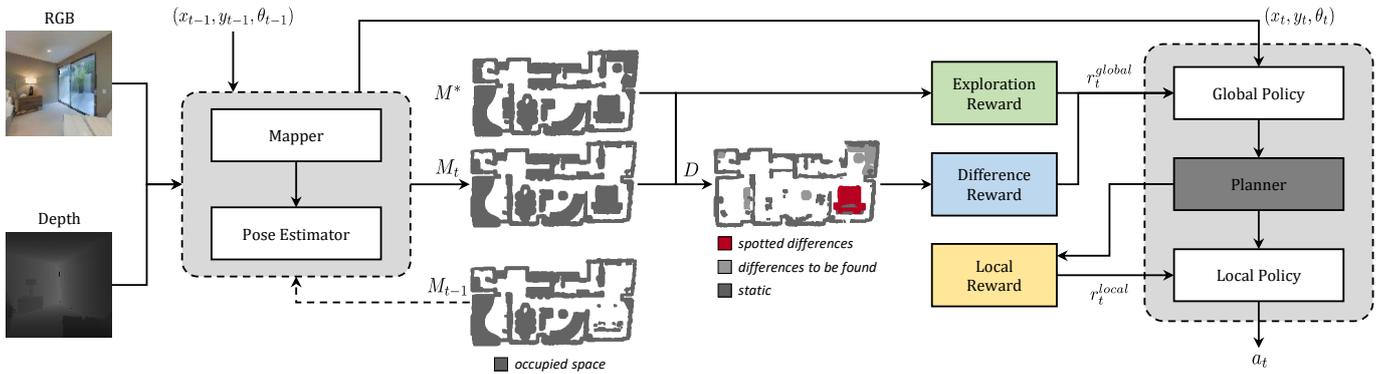


Fig. 2. Overview of the proposed approach for navigation in changing environments.

$\theta = 0^\circ$  (facing eastward). We assume that the input consists of a depth image and a semantic image and that the camera intrinsics  $K$  are known. To build the Semantic Occupancy Map (SOM) of an environment, we project each semantic pixel of the acquired scene into a 2-dimensional top-down map: given a pixel with image coordinates  $(i, j)$  and depth value  $d_{i,j}$ , we first recover its coordinates  $(x, y, z)$  with respect to the agent position. Then, we compute the corresponding  $(u, v)$  pixel in map through an orthographic projection, using the information about the agent position and heading:

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} = d_{i,j} K^{-1} \begin{bmatrix} i \\ j \\ 1 \end{bmatrix}, \quad \text{and} \quad \begin{bmatrix} u \\ v \\ 0 \\ 1 \end{bmatrix} = P_v \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}. \quad (2)$$

We perform the same operation after rotating the agent by  $\Delta\theta = 30^\circ$  until we perform a span from  $0^\circ$  to  $180^\circ$ . To cover the whole scene, we repeat this procedure placing the agent at a distance of  $0.5m$  from the previous capture point, following the axis directions. The agent elevation is instead kept fixed. During this step, we average the results of subsequent observations of overlapping portions of space.

After the acquisition, we obtain a SOM with  $C$  channels, where each pixel corresponds to a  $5cm \times 5cm$  portion of space in the 3D environment. For each channel  $c \in \{0, \dots, C\}$ , the map values represent the probability that the corresponding portion of space is occupied by an object of semantic class  $c$ . **Multiple Semantic Occupancy Maps for the Same Environment.** The SOMs obtained in the previous step can be seen as one possible layout for the corresponding 3D environments. In order to create a dataset with different states (*i.e.* different layouts) of the same environment, instead of manipulating the real-world 3D scenes (changing the furniture position, removing chairs, *etc.*), we propose to modify the SOM to create a set of plausible and different layouts for the environment.

First, we isolate the objects belonging to each semantic category by using an algorithm for connected component labeling [31], [32], [33]. Then, we sample a subset of objects to be deleted from the map and a subset of objects to be re-positioned in a different free location of the map. During

sampling, we consider categories that have a high probability of being displaced or removed in the real world and ignore non-movable semantic categories such as *fireplaces*, *columns*, and *stairs*. After this step, we obtain a new SOM representing a possible alternative state for the environment, which could be very different from the one in which the 3D acquisition was taken. Sample manipulated maps can be found in Fig. 1. **Dataset Details.** To generate alternative SOMs, we start from the Matterport 3D (MP3D) dataset of spaces [16], which comprises 90 different building scans, and is enriched with dense semantic annotations. We consider each floor in the building and compute the SOM for that floor. For each map, we create 10 alternative versions of that same environment. In this step, we discard the floors that have few semantic objects (*e.g.*, empty rooftops) or that are not fully navigable by the agent. As a result, we retain 249 floors belonging to 81 different buildings, thus generating a total of 2490 different semantic occupancy maps for these floors. Finally, we split the dataset into train, validation, and test subsets.

As an additional test bed, we also build a set of out-of-domain maps (13 floors from 5 spaces) taken from the Gibson dataset [17], enriched with semantic annotations from [34], and manipulated as done for the MP3D dataset. For each SOM, multiple episodes are generated by selecting different starting points. More information about our dataset can be found in Table I and in the supplementary material.

### C. Agent Architecture

Our model for embodied navigation in changing environments comprises three major components: a mapper module, a pose estimator, and a navigation policy (which, in turn, consists of a global policy, a planner, and a local policy). An overview of the proposed architecture is shown in Fig. 2 and described below, while additional details can be found in the supplementary material. Although the data we provide is enriched with semantic labels, our agent does not make use of such information directly. This is in line with current state-of-the-art architectures for embodied exploration that we choose as competitors.

**Mapper.** The mapper module takes as inputs an RGB observation  $o_t^r$  and the corresponding depth image  $o_t^d$ , representing

the first-person view of the agent at time-step  $t$ , and outputs the agent-centric occupancy map  $v_t$  of a  $V \times V$  region in front of the camera. Each pixel in  $v_t$  corresponds to a  $25mm \times 25mm$  portion of space and consists of two channels containing the probability of that cell being occupied and explored, respectively. As a first step, we encode  $o_t^r$  using the first two blocks of ResNet-18 pre-trained on ImageNet, followed by a three-layer CNN. We project the depth image  $o_t^d$  using the camera intrinsics [12] and obtain a preliminary map for the visible occupancy. We name the obtained feature representations  $\hat{o}_t^r$  and  $\hat{o}_t^d$ , respectively. We then encode the two feature maps using a U-Net [35]:

$$f_\mu(\hat{o}_t^r, \hat{o}_t^d) = \text{U-Net}_{\text{enc}}(\hat{o}_t^r, \hat{o}_t^d, \mu), \quad (3)$$

and decode the  $2 \times V \times V$  matrix of probabilities as:

$$v_t = \sigma(\text{U-Net}_{\text{dec}}(f_\mu(\hat{o}_t^r, \hat{o}_t^d), \phi)), \quad (4)$$

where  $\mu$  and  $\phi$  represent the learnable parameters in the U-Net encoder and decoder, respectively, and  $\sigma$  is the sigmoid activation function. The computed agent-centric occupancy map  $v_t$  is then registered in the global occupancy map  $M_{t-1}$  coming from the previous time-step to obtain  $M_t$ . To that end, we use a geometric transformation to project  $v_t$  in the global coordinate system, for which we need a triple  $(x, y, \theta)$  corresponding to the agent position and heading in the environment. This triple is estimated by a specific component that tracks the agent displacements across the environment, as discussed in the following paragraph.

**Pose Estimator.** The agent can move across the environment using three actions: *go forward 0.25m*, *turn left 10°*, *turn right 10°*. Since each action may produce a different outcome because of physical interactions with the environment (e.g., bumping into a wall) or noise in the actuation system, the pose estimator is used to estimate the real displacement made at every time-step. We estimate the agent displacement  $(\Delta x_t, \Delta y_t, \Delta \theta_t)$  at time-step  $t$  by using two consecutive RGB and depth observations, as well as the agent-centric occupancy maps  $(v_{t-1}, v_t)$  computed by the mapper at  $t-1$  and  $t$ . The actual agent position  $(x_t, y_t, \theta_t)$  is computed iteratively as:

$$(x_t, y_t, \theta_t) = (x_{t-1}, y_{t-1}, \theta_{t-1}) + (\Delta x_t, \Delta y_t, \Delta \theta_t). \quad (5)$$

We assume that the agent starting position is the triple  $(x_0, y_0, \theta_0) = (0, 0, 0)$ .

**Global Policy, Planner, and Local Policy.** The sampling of atomic actions for the exploration relies on a three-component hierarchical policy. The first component is the global policy, which samples a long-term global goal on the map. The global policy outputs a probability distribution over discretized locations of the global map. We sample the global goal from this distribution and then transform it in  $(x, y)$  global coordinates. The second component is a planner module, which employs the A\* algorithm to decode a local goal on the map. The local goal is an intermediate point, within  $0.25m$  from the agent, along the trajectory towards the global goal. The last element of our navigation module is the local policy, which decodes

the series of atomic actions taking the agent towards the local goal. In particular, the local policy is an RNN decoding the atomic action  $a_t$  to execute at every time-step. The reward  $r_t^{\text{local}}$  given to the local policy is proportional to the reduction in the Euclidean distance  $d$  between the agent position and the current local goal:

$$r_t^{\text{local}} = d_t - d_{t-1}. \quad (6)$$

Following the hierarchical structure, a global goal is sampled every  $N$  time-steps. A new local goal is computed if a new global goal is sampled, if the previous local goal is reached, or if the local goal location is known to be not traversable.

**Exploiting Past Knowledge for Efficient Navigation.** The global policy is trained using a two-term reward. The first term encourages exhaustive exploration and is proportional either to the increase of area-coverage [12] or to the increase of anticipated map accuracy as in [4]. Intuitively, the agent strives to maximize the portion of the seen area and thus maximizes the knowledge gathered during exploration. Moreover, since we consider a setting where a significant amount of knowledge is already available to the agent, we add a reward term to guide the agent towards meaningful points of the map. These correspond to the coordinates where major changes are likely to happen.

Given the occupancy map of the agent at time  $t$ ,  $M_t$ , the true occupancy map for the same environment  $M^*$ , and a time budget of  $T$  time-steps for exploration, we aim to minimize the following, for  $0 < t \leq T$ :

$$D = \sum \mathbb{1}[M_t \neq M^*] \quad (7)$$

In other words, we want to maximize the number of pixels in the online reconstructed map  $M_t$  that the agent correctly shifts from free to occupied (and vice-versa) during exploration. This leads to the reward term for difference discovery:

$$r_{\text{diff}} = \sum \mathbb{1}[M_t = M^*] - \sum \mathbb{1}[M_{t-1} = M^*]. \quad (8)$$

The proposed reward term is designed to encourage navigation towards areas in the map that are more likely to contain meaningful differences (e.g., rooms containing more objects that can be displaced or removed from the scene). Additionally, an agent trained with this reward will tend to avoid difficult spots that are likely to produce a mismatch in terms of the predicted occupancy maps. This is because errors in the mapping phase would result in a negative reward.

To train our model, we combine a reward promoting exploration and the more specific reward on found differences to exploit semantic clues in the environment:

$$r_t^{\text{global}} = \beta_1 r_{\text{exp}} + \beta_2 r_{\text{diff}} \quad (9)$$

where  $r_{\text{exp}}$  is the reward term encouraging task-agnostic exploration (such as coverage-based or anticipation-based rewards, as described in the next section), and  $\beta_1$  and  $\beta_2$  are two coefficients weighing the importance of the two elements.

TABLE II  
EXPERIMENTAL RESULTS ON MP3D TEST SET. THE AGENT INCORPORATING THE PROPOSED REWARD TERM FOR DISCOVERED DIFFERENCES OUTPERFORMS THE COMPETITORS ON THE MAIN METRICS FOR THE NOVEL SPOT THE DIFFERENCE TASK.

	Estimated Localization									Oracle Localization								
	Seen[%]	Acc.	IoU <sub>+</sub>	IoU <sub>-</sub>	IoU	mAcc.	mIoU <sub>+</sub>	mIoU <sub>-</sub>	mIoU	Seen[%]	Acc.	IoU <sub>+</sub>	IoU <sub>-</sub>	IoU	mAcc.	mIoU <sub>+</sub>	mIoU <sub>-</sub>	mIoU
<b>OccAnt</b>	52.1	26.2	13.4	6.1	11.5	51.1	19.1	8.3	15.8	49.0	35.6	26.5	16.1	24.8	77.8	49.2	23.6	43.2
<b>DR</b>	49.4	29.3	15.3	8.7	13.9	59.7	23.1	11.9	20.2	48.6	37.4	27.2	18.4	26.5	80.1	49.8	27.4	45.8
<b>AR</b>	43.8	30.6	19.7	12.9	18.8	72.5	36.8	18.4	32.7	43.6	32.5	23.2	17.5	23.0	78.7	47.5	26.7	44.5
<b>CR</b>	<b>53.2</b>	33.1	18.1	9.6	16.1	65.2	26.4	12.7	22.6	<b>52.8</b>	39.2	29.6	18.8	28.0	78.5	51.0	26.6	45.7
<b>AR+DR</b>	51.4	34.5	20.9	12.0	19.3	71.5	33.9	16.2	30.0	51.4	37.8	27.3	18.0	26.2	79.3	48.9	25.8	44.4
<b>CR+DR</b>	52.3	<b>37.8</b>	<b>24.2</b>	<b>14.8</b>	<b>22.7</b>	<b>76.2</b>	<b>39.1</b>	<b>19.8</b>	<b>34.8</b>	51.8	<b>40.3</b>	<b>29.2</b>	<b>19.2</b>	<b>28.1</b>	<b>82.1</b>	<b>50.4</b>	<b>26.9</b>	<b>46.2</b>

#### IV. EXPERIMENTS AND RESULTS

In this section, we detail our experimental setting and show experimental results for our new proposed task. Further analysis can be found in the supplementary material.

**Evaluation Metrics.** To evaluate the performance in *Spot the Difference*, we consider three main classes of metrics. First, we consider the percentage of navigable area in the environment seen by the agent during the episode (Seen[%]). Then, we evaluate the percentage of elements that have been correctly detected as changed in the occupancy map (Acc.) and the pixel-wise Intersection over Union for the *changed* occupancy map elements (IoU). Besides, we evaluate the task as a two-class problem and compute the IoU score for objects that were added in place of free space (IoU<sub>+</sub>) and for objects that were deleted during the map creation (IoU<sub>-</sub>). In addition, to evaluate the performance independently from the exploration capability, we propose to compute the metrics only on the portion of space that the agent actually visited (mAcc., mIoU, mIoU<sub>+</sub>, and mIoU<sub>-</sub>).

**Implementation Details.** We conduct our experiment using Habitat [18], a popular platform for Embodied AI in photo-realistic indoor environments [17], [16]. The agent observations are 128 × 128 RGB-D images from the environment. The learning algorithm adopted for training is PPO [36]. The learning rate is 10<sup>-3</sup> for the mapper and 2.5 × 10<sup>-4</sup> for the other modules. Every model is trained for ≈ 6.5M frames using Adam optimizer [37]. A global goal is sampled every N = 25 time-steps. The local and global policies are updated, respectively, every N and 20 × N time-steps, and the mapper is updated every 4 × N time-steps. The size of the local map is V = 101, while the global map size is set to W = 2001 for the MP3D dataset and to W = 961 for the Gibson dataset. The global policy action space size G is 240. The reward coefficients {β<sub>1</sub>, β<sub>2</sub>} are set to {1, 10<sup>-2</sup>} and {1, 10<sup>-1</sup>} when the exploration reward is based on coverage and anticipation reward, respectively. The length of each episode is fixed to T = 1000 time-steps.

**Competitors and Baselines.** We consider the following competitors and variants of the proposed method on two different setups: one where the agent position is predicted by the agent (as in Eq. 5), and one where it has access to oracle coordinates: *Difference Reward (DR)*: an exploration policy that maximizes the correctly predicted changes between M and M\*. This

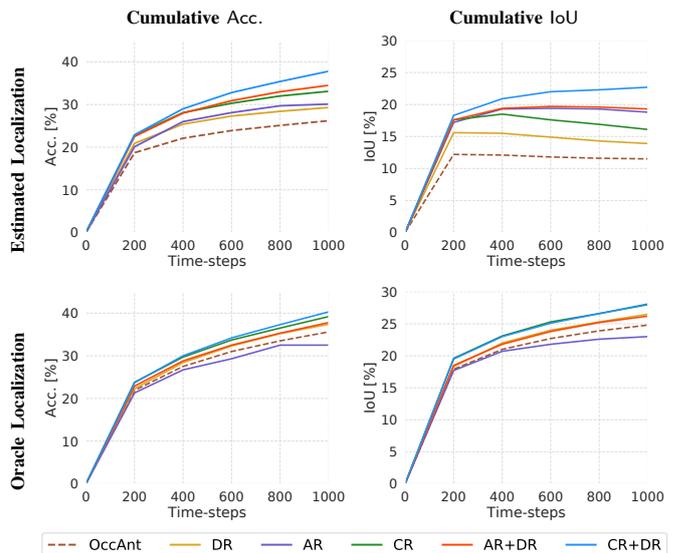


Fig. 3. Value of accuracy and IoU for the different models at varying time-steps on the MP3D test set.

corresponds to setting β<sub>1</sub> = 0 and β<sub>2</sub> = 1 in Eq. 9.

*Coverage Reward (CR)*: an agent that explores the environment with an exploration policy that maximizes the covered area and builds the occupancy map as it goes, as in [4].

*Anticipation Reward (AR)*: an agent that explores the environment with an exploration policy that maximizes the covered area and the correctly anticipated values in the occupancy map built as it goes, from [4]. Our proposed approach consists of an agent trained with the combination of the difference reward with the coverage reward (CR+DR) or with the anticipation reward (AR+DR).

*Occupancy Anticipation (OccAnt)*: we also compare with the agent presented by Ramakrishnan *et al.* [4] using the available pre-trained models, referenced to as *OccAnt*. Note that *OccAnt* was trained on the Gibson dataset for the standard exploration task and without any prior map. Thus, it is not directly comparable with the other methods considered. We include it to gain insights into the performance of an off-the-shelf state-of-the-art agent on our task.

**Results on MP3D dataset.** As a first testbed, we evaluate the different agents on the MP3D *Spot the Difference* test set. We report the results for this experiment in Table II.

We observe that the agent combining a reward based on

TABLE III  
EXPERIMENTAL RESULTS ON GIBSON VALIDATION SET.

	Estimated Localization								
	Seen[%]	Acc.	IoU <sub>+</sub>	IoU <sub>-</sub>	IoU	mAcc.	mIoU <sub>+</sub>	mIoU <sub>-</sub>	mIoU
<b>OccAnt</b>	86.2	49.8	11.9	7.2	10.4	58.0	12.3	7.5	10.8
<b>DR</b>	<b>86.2</b>	53.2	13.2	8.5	11.7	63.7	13.9	8.8	12.3
<b>AR</b>	75.3	51.5	21.4	16.6	20.4	72.7	25.8	17.3	23.3
<b>CR</b>	85.9	57.6	16.7	11.9	15.4	71.3	18.6	12.3	16.7
<b>AR+DR</b>	83.4	58.7	20.0	14.9	19.0	75.8	23.0	15.6	21.1
<b>CR+DR</b>	82.1	<b>60.1</b>	<b>24.0</b>	<b>19.0</b>	<b>23.1</b>	<b>78.5</b>	<b>27.8</b>	<b>19.9</b>	<b>25.9</b>

coverage and our reward based on the differences in the environment (*CR+DR*) performs best on all the pixel-based metrics and places second in terms of percentage of seen area. It is worth noting that, even if the results in terms of the area seen are not as high as the ones obtained by the *CR* agent, the addition of our Difference Reward helps the agent to focus on more relevant parts, and thus, it can discover more substantial differences. Additionally, predictions are more accurate and more precise, as indicated by the 4.7% and 6.6% improvements in terms of Acc. and IoU with respect to the *CR* competitor. Instead, a reward based on differences alone is not sufficient to promote good exploration. In fact, although the *DR* agent outperforms the *CR* and *AR* agents on some metrics, our reward alone does not provide as much improvement as when combined with rewards encouraging exploration (as for *CR+DR* and *AR+DR*).

Even in the oracle localization setup, the *CR+DR* agent achieves the best results. Interestingly, the gap with the *CR* agent decreases to 1.1% and 0.1% in terms of Acc. and IoU, respectively. This is because our *CR+DR* agent learns to sample trajectories that can be performed more efficiently and without accumulating a high positioning error. For this reason, the performance boost given by the oracle localization is lower. For both setups, our *CR+DR* agent outperforms the state-of-the-art *OccAnt* agent for exploration on all the metrics.

Finally, in Fig. 3, we plot different values of Acc. and IoU over different time-steps during the episodes. This way, we can evaluate the whole exploration trend, and not only its final point. We can observe that the proposed models incorporating the difference reward outperform the competitors. In particular, the *CR+DR* agent scores first by a significant margin. The performance gap can be noticed even in the first half of the episode and tends to grow with the number of steps.

**Results on Gibson dataset.** The environments from the Gibson dataset [17] are generally smaller than those in MP3D, and thus, they can be explored more easily and exhaustively. We report the results for this experiment in Table III. Also in this experiment, the *CR+DR* agent performs best on all the metrics but the percentage of the area seen. Although *CR+DR* explores 3.8% of the environment less than the *CR* agent, it still overcomes the competitor by 2.5% and 7.7% in terms of Acc. and IoU. The *AR+DR* agent is the second-best in terms of Acc.. The *OccAnt* agent, instead, is competitive in terms of area seen but achieves low Acc. and IoU metrics.

**Qualitative Results.** In Fig. 4, we report some qualitative



Fig. 4. Qualitative results comparing the performances of the *CR* and *CR+DR* agents for different episodes.

results. Starting from the left-most column, we present the starting map given to the agent as the episode begins, the results achieved by the *CR* agent, those of the proposed *CR+DR* agent, and the ground-truth map. The differences that the agents have correctly identified during the episode are highlighted in red. As it can be seen, the *CR+DR* agent can identify more differences than the *CR* counterpart, even in small environments (top row). As the size of the environments grows (bottom row), the performance gap increases and the *CR+DR* agent outperforms its competitor.

## V. CONCLUSION

In this work, we proposed *Spot the Difference*: a new task for navigation agents in changing environments. In this novel setting, the agent has to find all variations that occurred in the environment with respect to an outdated occupancy map. Since current datasets of 3D spaces do not account for such variety, we collected a new dataset containing different layouts for the same environment. We tested two state-of-the-art exploration agents on this task and proposed a novel reward term to encourage the discovery of meaningful information during exploration. The proposed agent outperforms the competitors and can identify changes in the environment more efficiently. We believe that the results presented in this paper motivate further research on this new proposed setting for Embodied AI.

## ACKNOWLEDGMENT

This work has been supported by the “European Training Network on Personalized Robotics as Service Oriented applications” (PERSEO) MSCA-ITN-2020 project (G.A. 955778).

## REFERENCES

- [1] S. K. Ramakrishnan, D. Jayaraman, and K. Grauman, "An Exploration of Embodied Visual Exploration," *arXiv preprint arXiv:2001.02192*, 2020.
- [2] R. Bigazzi, F. Landi, M. Cornia, S. Cascianelli, L. Baraldi, and R. Cucchiara, "Explore and Explain: Self-supervised Navigation and Recounting," in *Proceedings of the International Conference on Pattern Recognition*, 2020.
- [3] R. Bigazzi, F. Landi, S. Cascianelli, L. Baraldi, M. Cornia, and R. Cucchiara, "Focus on Impact: Indoor Exploration with Intrinsic Motivation," *IEEE Robotics and Automation Letters*, 2022.
- [4] S. K. Ramakrishnan, Z. Al-Halah, and K. Grauman, "Occupancy Anticipation for Efficient Exploration and Navigation," in *Proceedings of the European Conference on Computer Vision*, 2020.
- [5] D. S. Chaplot, D. P. Gandhi, A. Gupta, and R. R. Salakhutdinov, "Object Goal Navigation using Goal-Oriented Semantic Exploration," in *Advances in Neural Information Processing Systems*, 2020.
- [6] J. Krantz, E. Wijmans, A. Majumdar, D. Batra, and S. Lee, "Beyond the nav-graph: Vision-and-language navigation in continuous environments," *Proceedings of the European Conference on Computer Vision*, 2020.
- [7] F. Landi, L. Baraldi, M. Cornia, M. Corsini, and R. Cucchiara, "Multimodal Attention Networks for Low-Level Vision-and-Language Navigation," *Computer Vision and Image Understanding*, vol. 210, p. 103255, 2021.
- [8] V. Cartillier, Z. Ren, N. Jain, S. Lee, I. Essa, and D. Batra, "Semantic MapNet: Building Allocentric Semantic Maps and Representations from Egocentric Views," *arXiv preprint arXiv:2010.01191*, 2020.
- [9] S. Wani, S. Patel, U. Jain, A. X. Chang, and M. Savva, "MultiON: Benchmarking Semantic Map Memory using Multi-Object Navigation," in *Advances in Neural Information Processing Systems*, 2020.
- [10] D. S. Chaplot, D. Gandhi, S. Gupta, A. Gupta, and R. Salakhutdinov, "Learning To Explore Using Active Neural SLAM," in *Proceedings of the International Conference on Learning Representations*, 2019.
- [11] P. Anderson, A. Chang, D. S. Chaplot, A. Dosovitskiy, S. Gupta, V. Koltun, J. Kosecka, J. Malik, R. Mottaghi, M. Savva *et al.*, "On evaluation of embodied navigation agents," *arXiv preprint arXiv:1807.06757*, 2018.
- [12] T. Chen, S. Gupta, and A. Gupta, "Learning Exploration Policies for Navigation," in *Proceedings of the International Conference on Learning Representations*, 2019.
- [13] M. Luperto, M. Antonazzi, F. Amigoni, and N. A. Borghese, "Robot exploration of indoor environments using incomplete and inaccurate prior knowledge," *Robotics and Autonomous Systems*, vol. 133, p. 103622, 2020.
- [14] P. Karkus, S. Cai, and D. Hsu, "Differentiable SLAM-net: Learning Particle SLAM for Visual Navigation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021.
- [15] B. Mayo, T. Hazan, and A. Tal, "Visual navigation with spatial attention," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021.
- [16] A. Chang, A. Dai, T. Funkhouser, M. Halber, M. Niessner, M. Savva, S. Song, A. Zeng, and Y. Zhang, "Matterport3D: Learning from RGB-D Data in Indoor Environments," in *Proceedings of the International Conference on 3D Vision*, 2017.
- [17] F. Xia, A. R. Zamir, Z. He, A. Sax, J. Malik, and S. Savarese, "Gibson Env: Real-world perception for embodied agents," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [18] M. Savva, A. Kadian, O. Maksymets, Y. Zhao, E. Wijmans, B. Jain, J. Straub, J. Liu, V. Koltun, J. Malik *et al.*, "Habitat: A Platform for Embodied AI Research," in *Proceedings of the International Conference on Computer Vision*, 2019.
- [19] A. Kadian, J. Truong, A. Gokaslan, A. Clegg, E. Wijmans, S. Lee, M. Savva, S. Chernova, and D. Batra, "Sim2Real Predictivity: Does evaluation in simulation predict real-world performance?" *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 6670–6677, 2020.
- [20] R. Bigazzi, F. Landi, M. Cornia, S. Cascianelli, L. Baraldi, and R. Cucchiara, "Out of the Box: Embodied Navigation in the Real World," in *Proceedings of the International Conference on Computer Analysis of Images and Patterns*, 2021.
- [21] S. Gupta, J. Davidson, S. Levine, R. Sukthankar, and J. Malik, "Cognitive mapping and planning for visual navigation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [22] F. L. Da Silva, M. E. Taylor, and A. H. R. Costa, "Autonomously Reusing Knowledge in Multiagent Reinforcement Learning," in *Proceedings of the International Joint Conferences on Artificial Intelligence*, 2018.
- [23] D. S. Chaplot, L. Lee, R. Salakhutdinov, D. Parikh, and D. Batra, "Embodied Multimodal Multitask Learning," *Proceedings of the International Joint Conferences on Artificial Intelligence*, 2019.
- [24] N. Savinov, A. Dosovitskiy, and V. Koltun, "Semi-parametric topological memory for navigation," in *Proceedings of the International Conference on Learning Representations*, 2018.
- [25] M. Sridharan and T. Mota, "Commonsense Reasoning to Guide Deep Learning for Scene Understanding," in *Proceedings of the International Joint Conferences on Artificial Intelligence*, 2020.
- [26] Y. Zhang, H. Tan, and M. Bansal, "Diagnosing the Environment Bias in Vision-and-Language Navigation," *Proceedings of the International Joint Conferences on Artificial Intelligence*, 2020.
- [27] M. R. U. Saputra, A. Markham, and N. Trigoni, "Visual SLAM and structure from motion in dynamic environments: A survey," *ACM Computing Surveys*, vol. 51, no. 2, pp. 1–36, 2018.
- [28] J. Biswas, "The Quest For" Always-On" Autonomous Mobile Robots," in *Proceedings of the International Joint Conferences on Artificial Intelligence*, 2019.
- [29] T. T. Mac, C. Copot, D. T. Tran, and R. De Keyser, "Heuristic approaches in robot path planning: A survey," *Robotics and Autonomous Systems*, vol. 86, pp. 13–28, 2016.
- [30] L. Nardi and C. Stachniss, "Long-term robot navigation in indoor environments estimating patterns in traversability changes," in *Proceedings of the IEEE International Conference on Robotics and Automation*, 2020.
- [31] C. Grana, D. Borghesani, and R. Cucchiara, "Optimized block-based connected components labeling with decision trees," *IEEE Transactions on Image Processing*, vol. 19, no. 6, pp. 1596–1609, 2010.
- [32] F. Bolelli, S. Allegretti, L. Baraldi, and C. Grana, "Spaghetti labeling: Directed acyclic graphs for block-based connected components labeling," *IEEE Transactions on Image Processing*, vol. 29, pp. 1999–2012, 2019.
- [33] S. Allegretti, F. Bolelli, and C. Grana, "Optimized block-based algorithms to label connected components on gpus," *IEEE Transactions on Parallel and Distributed Systems*, vol. 31, no. 2, pp. 423–438, 2019.
- [34] I. Armeni, Z.-Y. He, J. Gwak, A. R. Zamir, M. Fischer, J. Malik, and S. Savarese, "3D Scene Graph: A structure for unified semantics, 3D space, and camera," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [35] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation," in *Proceedings of the International Conference on Medical Image Computing and Computer Assisted Intervention*, 2015.
- [36] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [37] D. Kingma and J. Ba, "Adam: a method for stochastic optimization," in *Proceedings of the International Conference on Learning Representations*, 2015.

### A. Additional Implementation Details

**Semantic Classes Division.** The generation of semantic maps for each floor of each scene produces  $2001 \times 2001 \times 43$  maps for the MP3D dataset and  $961 \times 961 \times 21$  maps for the Gibson dataset. The last channel of every map registers the explorable space, so it is ignored for the creation of the dataset and is concatenated, as it is, to the manipulated map obtained at the end of the semi-automatic dataset creation process.

We divide the semantic channels of the maps depending on the possible actions performable on the connected components in that channel. We identify four types of classes: *No Operation*, *Removal*, *Displacement*, and *Overlap Removal*. A list of semantic categories with their classification is reported in Table VI for the MP3D dataset and in Table VII for the Gibson dataset. *No Operation* classes are left untouched, and correspond to non movable objects, such as *wall*, *stairs* and *columns*; the connected component of the *Removal* classes can be removed; those in the *Displacement* classes can be either removed or relocated in other free spaces in the environment; and *Overlap Removal* components are removed if connected components removed or displaced in other channels overlap with them, e.g., if a *sofa* is removed, every instance of *cushion* overlapping with that *sofa* will be removed as well because it is supposed to be on top of it.

In Fig. 9, we report some additional examples of manipulated semantic maps with relative difference maps obtained by applying our semi-automatic procedure.

**Train/Val/Test Splits Division.** We use the same scene partitioning adopted by the existing datasets for embodied exploration and PointGoal navigation on Matterport3D and Gibson Tiny [16], [17].

**Episode Creation.** For the creation of the episodes of our dataset we use the starting positions of the exploration dataset for MP3D, and of the PointGoal navigation dataset for Gibson Tiny. After the episodes located in floors with few semantic objects or that are not fully navigable by the agent are discarded, we associate one of the alternative versions of the ground-truth semantic map to each episode. For the validation and test splits of the MP3D dataset and the validation split of the Gibson dataset we create new episodes with random sampled starting positions so that the number of episodes on every floor is at least 10 and fix the number of episodes per floor to a multiple of 10. We report a detailed list of scans, selected floors and number of episodes per scan in Tables IX, X, XI, and XII.

**Pose Estimator.** The pose estimator takes as input two consecutive RGB and depth observations, consisting in two pairs  $(o_{t-1}^r, o_{t-1}^d)$  and  $(o_t^r, o_t^d)$ . Additionally, it accepts as input the agent-centric occupancy maps  $(v_{t-1}, v_t)$  computed by the mapper at  $t-1$  and  $t$ . For each modality, we encode information using a CNN followed by a fully-connected layer. We call these intermediate representations  $\bar{o}_t^r$ ,  $\bar{o}_t^d$ , and  $\bar{v}_t$ . Then,

we compute a first estimate of the relative displacement in terms of  $(x, y, \theta)$  coordinates and heading for each modality:

$$g(\star) = W_1 \max(W_2 \star + b_2, 0) + b_1, \quad (10)$$

with  $\star \in \{\bar{o}_t^r, \bar{o}_t^d, \bar{v}_t\}$ . We stack the vectors computed in Eq. 10 to obtain a  $3 \times 3$  matrix  $G$ . Finally, we compute the agent displacement at time-step  $t$   $(\Delta x_t, \Delta y_t, \Delta \theta_t)$  as:

$$(\Delta x_t, \Delta y_t, \Delta \theta_t) = \sum_{i=1}^3 \alpha_i \cdot G_i, \quad (11)$$

$$\alpha_i = \text{softmax}(\text{MLP}_i([\bar{o}_t^r, \bar{o}_t^d, \bar{v}_t])), \quad (12)$$

where  $G_i$  indicates the  $i$ -th row of the  $G$  matrix, MLP is a three-layer fully-connected network, and  $[\cdot, \cdot, \cdot]$  denotes tensor concatenation.

**Global Policy.** An enriched occupancy map  $M_t^+ \in [0, 1]^{4 \times W \times W}$  is obtained by stacking the occupancy map, the map of visited states, and the one-hot representation of the agent location  $(x_t, y_t)$ . Then, we compute two versions of  $M_t^+$ : one by cropping the map to an agent-centered  $G \times G$  area, and the other by max-pooling the map to the same spatial resolution. The 8-channel tensor obtained by concatenating these two versions of  $M_t^+$  is fed to the global policy. The global policy consists of a CNN that outputs a probability distribution over the  $G \times G$  global action space. We sample the global goal from this distribution, and then transform it in  $(x, y)$  global coordinates.

**List of Hyperparameters.** In Table VIII, we list the hyperparameters used to train our agents. These hyperparameters are shared among all the agents. Agent-depending values are specified in the main paper.

**Training setup.** We train every agent using 16 NVIDIA V100 GPUs in parallel, distributed on 4 different nodes with 4 GPUs each. On every node, we run 8 Habitat environments, for a total of 32 environments in parallel. To coordinate the interaction among nodes and train our models, we use PyTorch. Each train took about 48 hours using this setup.

### B. Additional Experimental Results

**Results on MP3D Validation Set.** We report the quantitative results on the MP3D validation set in Table IV. Discussion for these experiments is analogous to the one presented for the experiments on the MP3D test set in the main paper.

**Results on Gibson Validation Set with Oracle Localization.** We report the quantitative results on the Gibson validation set using the oracle localization setup in Table V. In this setting, the agent using only the proposed difference reward (*DR*) performs the best on almost all the metrics. We can conclude that, for small environments, and given an optimal localization system, our reward alone is sufficient to surpass the competitors on *Spot the Difference*.

**Analysis at Different Time-Steps.** In Fig. 5 and Fig. 6, we report the plots of different values of Acc. and IoU for the MP3D and Gibson validation sets, respectively. In these plots, we show how Acc. and IoU vary at different time-steps during the episodes for the various methods.

TABLE IV  
EXPERIMENTAL RESULTS ON MP3D VALIDATION SET.

	Estimated Localization									Oracle Localization								
	Seen[%]	Acc.	IoU <sub>+</sub>	IoU <sub>-</sub>	IoU	mAcc.	mIoU <sub>+</sub>	mIoU <sub>-</sub>	mIoU	Seen[%]	Acc.	IoU <sub>+</sub>	IoU <sub>-</sub>	IoU	mAcc.	mIoU <sub>+</sub>	mIoU <sub>-</sub>	mIoU
<b>OccAnt</b> <sup>†</sup>	50.3	24.1	9.6	5.8	8.5	49.5	13.9	7.0	11.6	47.5	33.5	21.5	17.4	20.6	76.0	41.6	23.0	20.1
<b>DR</b>	48.0	28.7	11.9	7.3	10.7	59.7	17.3	9.0	14.7	44.7	35.2	22.3	19.0	22.1	79.4	42.4	25.0	38.2
<b>AR</b>	40.1	29.4	15.9	13.0	15.7	71.3	28.8	<b>16.8</b>	25.9	39.2	31.0	19.2	17.7	19.4	77.6	38.9	24.8	36.0
<b>CR</b>	<b>53.2</b>	34.7	14.6	8.7	12.9	64.1	20.3	10.3	17.0	<b>52.2</b>	<b>42.6</b>	<b>26.9</b>	19.8	<b>25.4</b>	80.4	<b>46.4</b>	25.2	<b>40.2</b>
<b>AR+DR</b>	48.7	34.9	16.6	11.9	15.7	71.7	26.7	14.2	22.9	47.6	38.0	22.6	19.3	22.4	<b>81.3</b>	42.0	25.4	37.5
<b>CR+DR</b>	50.1	<b>38.6</b>	<b>20.1</b>	<b>14.1</b>	<b>19.1</b>	<b>75.4</b>	<b>31.8</b>	<b>16.8</b>	<b>27.6</b>	50.0	40.9	25.1	<b>19.9</b>	24.3	80.8	43.8	<b>26.2</b>	39.2

TABLE V  
EXPERIMENTAL RESULTS ON GIBSON VALIDATION SET USING THE ORACLE LOCALIZATION SETUP.

	Oracle Localization								
	Seen[%]	Acc.	IoU <sub>+</sub>	IoU <sub>-</sub>	IoU	mAcc.	mIoU <sub>+</sub>	mIoU <sub>-</sub>	mIoU
<b>OccAnt</b>	81.6	60.1	<b>32.1</b>	21.2	29.2	78.7	<b>39.6</b>	22.2	<b>34.1</b>
<b>DR</b>	<b>86.1</b>	<b>65.2</b>	30.1	<b>24.1</b>	<b>29.9</b>	81.1	36.0	25.2	33.8
<b>AR</b>	74.1	53.8	27.9	21.9	27.2	77.0	35.4	23.5	32.7
<b>CR</b>	84.0	62.2	30.6	22.1	28.8	79.5	36.1	23.3	32.8
<b>AR+DR</b>	83.2	63.2	29.6	23.8	29.1	<b>81.6</b>	35.8	25.1	33.7
<b>CR+DR</b>	82.6	63.8	30.3	<b>24.1</b>	29.5	<b>81.6</b>	36.1	<b>25.5</b>	34.0

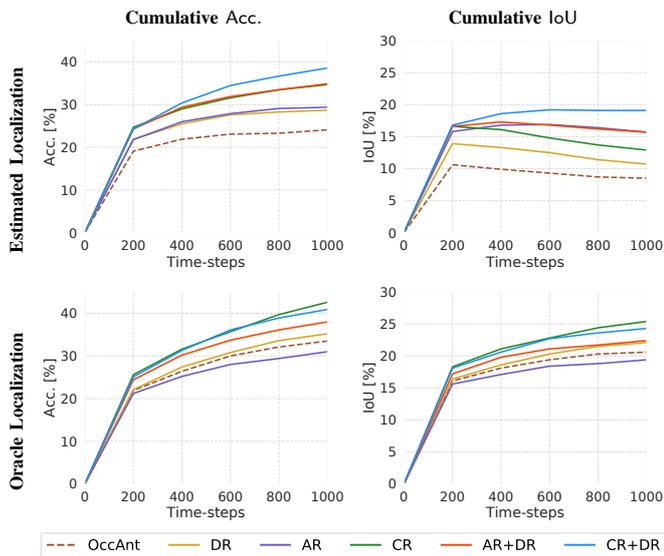


Fig. 5. Value of accuracy and IoU for the different models at varying time-steps on the MP3D validation set.

**Qualitative Results.** We include some additional qualitative results for the proposed *Spot the Difference* task in Fig. 8. In particular, we compare the CR agent with the CR+DR counterpart on different episodes with different map size and complexity. For each episode, we report the starting map given to the agent, the reconstructed map collected after exploration, and the ground-truth map of the actual state of the environment. We display the spotted differences in red while keeping the undiscovered elements in light gray. Dark gray denotes unchanged elements of the environment. The proposed CR+DR agent can discover a larger set of differences, thus

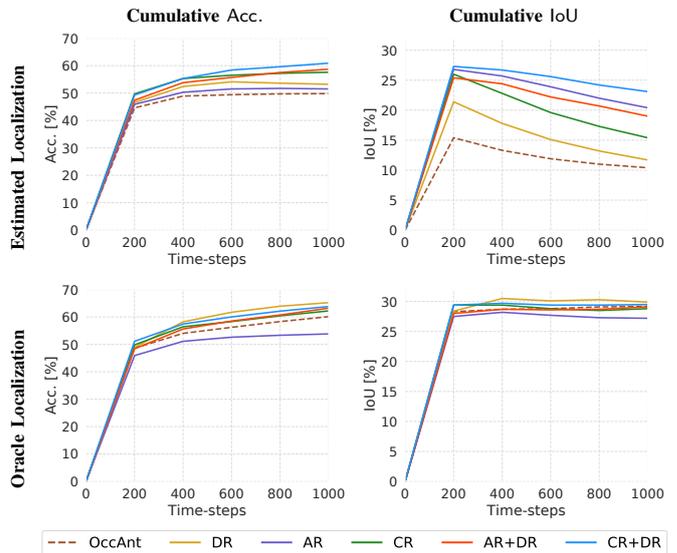


Fig. 6. Value of accuracy and IoU for the different models at varying time-steps on the Gibson validation set.

achieving better results.

Moreover, we report sample exploration trajectories for the CR and the CR+DR agents with estimated localization in Fig. 7. These confirm the competitive exploration capabilities of our proposed agent.

### C. Discussion and Future Directions

We present a method that exploits outdated information about the current environment to improve the exploration capabilities of the agent. However, the focus of this work is on pure occupation, ignoring semantic information. For future work, we expect to include semantic reasoning into the agent’s pipeline. We assume that additional information could boost the performance. With the proposed dataset, we enable a series of possible embodied tasks that imply dynamic environments and incorporate available past knowledge.

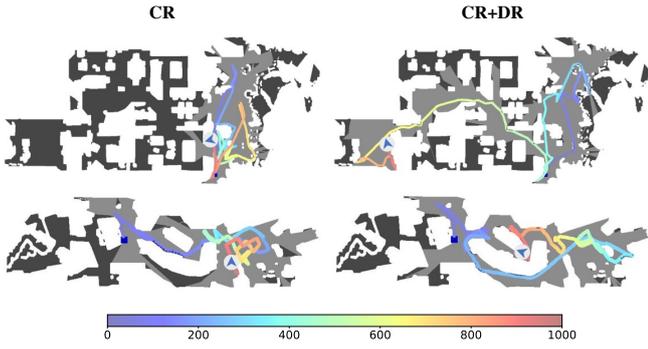


Fig. 7. Exploration trajectories of the CR and CR+DR agents on sample MP3D test episodes.



Fig. 8. Qualitative results comparing the performances of the CR and CR+DR agents for different episodes.

TABLE VI  
MP3D SEMANTIC CATEGORIES PER CHANNEL INDEX.

MP3D		
Index	Category	Action
0	Void	No Operation
1	Wall	No Operation
2	Floor	No Operation
3	Chair	Displacement
4	Door	No Operation
5	Table	Displacement
6	Picture	No Operation
7	Cabinet	Removal
8	Cushion	Overlap Removal
9	Window	No Operation
10	Sofa	Displacement
11	Bed	Displacement
12	Curtain	No Operation
13	Chest of Drawers	Displacement
14	Plant	Displacement
15	Sink	Empty
16	Stairs	No Operation
17	Ceiling	No Operation
18	Toilet	Removal
19	Stool	Displacement
20	Towel	Overlap Removal
21	Mirror	No Operation
22	TV Monitor	Removal
23	Shower	Removal
24	Column	No Operation
25	Bathtub	Removal
26	Counter	Removal
27	Fireplace	No Operation
28	Lighting	No Operation
29	Beam	No Operation
30	Railing	No Operation
31	Shelving	Removal
32	Blinds	No Operation
33	Gym Equipment	Displacement
34	Seating	Removal
35	Board Panel	No Operation
36	Furniture	Displacement
37	Appliances	Removal
38	Clothes	Overlap Removal
39	Objects	Overlap Removal
40	Misc	Overlap Removal
41	Unlabeled	No Operation

TABLE VII  
GIBSON SEMANTIC CATEGORIES PER CHANNEL INDEX.

Gibson		
Index	Category	Action
0	Chair	Displacement
1	Couch	Displacement
2	Potted Plant	Removal
3	Bed	Displacement
4	Toilet	Removal
5	TV	Removal
6	Dining Table	Displacement
7	Oven	Removal
8	Sink	Removal
9	Refrigerator	Removal
10	Book	Overlap Removal
11	Clock	Removal
12	Vase	Removal
13	Cup	Overlap Removal
14	Bottle	Overlap Removal
15	Bench	Removal
16	Appliances	Removal
17	Objects	Overlap Removal
18	Misc	Overlap Removal
19	Void	No Operation

TABLE VIII  
LIST OF HYPERPARAMETERS.

Hyperparameters		
Module	Name	Value
Global Policy	map size (G)	240
	lr	$2.5 \times 10^{-4}$
	max_grad_norm	0.5
Local Policy	forward step	$0.25m$
	turn angle	$10^\circ$
	hidden size	256
	lr	$2.5 \times 10^{-4}$
Mapper	batch size	32
	map scale	$0.05cm^2$
	map size	101
	lr	$10^{-3}$
	momentum	0.9
	max_grad_norm	0.5
PPO	clip param	0.2
	entropy coef	$10^{-3}$
	eps	$10^{-5}$
	gamma	0.99
	n. mini-batches	4
	n. epochs	4
	tau	0.95
	using gae	True
	value loss coef	0.5

TABLE IX  
MP3D TRAIN SCANS AND FLOORS.

MP3D Train		
Scan	Floors	# Episodes
HxpKQynjfin	0	81967
gTV8FGcVJC9	0,1,2,3,4,6,10,11	77186
29hnd4uzFmX	0,1,2,3	81967
5LpN3gDmAk7	0,1,2,3	81885
SN83YJsR3w2	0,1,2,3,7,8,10,12	81438
VzqfbhrpDEA	0,1,3,6	81641
D7N2EKCX4Sj	0,1,2,3,5,6	81830
5q7pvUzZiYa	0,1,2,3,4	81967
ac26ZMwG7aT	0,1	81967
r47D5H71a5s	0,1	81965
Pm6F8kyY3z2	0	81967
8WUmhLawc2A	0,1,2	81967
82sE5b5pLXE	0,1,2	80682
mJXqzFtmKg4	0,1,2	81967
i5noydFURQK	0,1	81120
V2XKFyX4ASd	0,1,2,3,4,5,7	81129
759xd9YjKW5	0,1,2,3	81913
r1Q1Z4BcV1o	0	81812
S9hNv5qa7GM	0,1	81967
1LXtFkfw3qL	0,1,2,3,4,5,6	81967
PuKPg4mmafe	0	81940
EDJbREhghzL	0,1,3	64755
ur6pFq6Qu1A	0,1	81967
B6ByNegPMKs	0	81951
b8cTxDM8gDG	0,1,2,8,11	73307
17DRP5sb8fy	0	81967
YmJkqBEsHnH	0	80780
ULsKaCPVFJR	0,1,2	81967
XcA2TqTSSAj	0,2,3,5,6,8,9,11,12	60679
sKLMlPTheUy	0,1,2,4	79736
ZMojNkEp431	0,1,2	81967
e9zR4mvMWw7	0,1,2	80193
JeFG25nYj2p	0,1	81967
uNb9QFRL6hY	1,4,5,6	59613
p5wJjkQkbXX	0,1,2,3	81967
Vvot9Ly1tCj	0,3	78115
E9uDoFAP3SH	0,1,5,6	81914
qoiz87JEwZ2	0,1,2,3	81967
VFuaQ6m2Qom	0,1,2,4,5,6	81758
VLzqgDo317F	0,1,2	81396
kEZ7cmS4wCh	0,1,2,3,7	69135
7y3sRwLe3Va	0,1,2,5	81386
VVfe2KiqLaN	0,1,2	81967
2n8kARJN3HM	0,1,2,4	81076
PX4nDJXEHRG	0,1,2,3,4,5	79151
Uxmj2M2itWa	0,1,3,4	49942
pRbA3pwrkg9	0,2,3,7,9,11	53295
cV4RvEzvu5T	0,1,2,3	81038
sT4fr6TAbpF	0	81625
GdvgFV5R1Z5	0	81967
JF19kD82Mey	0,1,2	81927
JmbYfDe2QKZ	0,1	81489
s8pcmisQ38h	0,1,2	80428
1pXnuDYAj8r	0,1,2,5	81901
jh4fe5c5qoQ	0,1,2	81967
vyrNrziPKCB	0,1,3,4,7	81388
aayBHfsNo7d	0,1,2	81693
rPc6DW4iMge	0,1,3,4	80296
<b>Total: 58</b>	<b>207</b>	<b>4581881</b>

TABLE X  
MP3D VALIDATION SCANS AND FLOORS, WITH RELATIVE NUMBER OF EPISODES FOR *Spot the Difference*.

MP3D Validation		
Scan	Floors	# Episodes
2azQ1b91cZZ	0,1	40
8194nk5LbLH	0	40
EU6Fwq7SyZv	0	30
QUCTc6BB5sX	1	20
TbHJrupSAjP	0,1,2	30
Z6MFQCvIBuw	0	40
oLBMNvg9in8	0,1,2,3	50
x8F5xyUWY9e	0,1	30
zsNo4HB9uLZ	0	40
<b>Total: 9</b>	<b>16</b>	<b>320</b>

TABLE XI  
MP3D TEST SCANS AND FLOORS, WITH RELATIVE NUMBER OF EPISODES FOR *Spot the Difference*.

MP3D Test		
Scan	Floors	# Episodes
2t7WUuJeko7	0	50
5ZKStnWn8Zo	0,1	50
RPmz2sHmrrY	0	50
UwV83HsGsw3	0,1,2,3	50
WYY7iVvf5p8	0,2	30
YFuZgdQ5vWj	1	10
YVUC4YcDtcY	0	50
fzynW3qQPVF	0,1	50
jtcxE69GiFV	0,1	40
pa4otMbVnkk	0,1	50
q9vSo1VnCiC	0	50
rqfALeAoiTq	0,2	20
wc2JMjhGNzB	0,1	50
yqstnuAEVhm	0,1,2	60
<b>Total: 14</b>	<b>26</b>	<b>610</b>

TABLE XII  
GIBSON VALIDATION SCANS AND FLOORS, WITH RELATIVE NUMBER OF EPISODES FOR *Spot the Difference*.

Gibson Validation		
Scan	Floors	# Episodes
Wiconisco	1,2	90
Corozal	0,2,4	90
Collierville	0,1,2	80
Markleeville	0,1	90
Darden	0,1,2	100
<b>Total: 5</b>	<b>13</b>	<b>450</b>

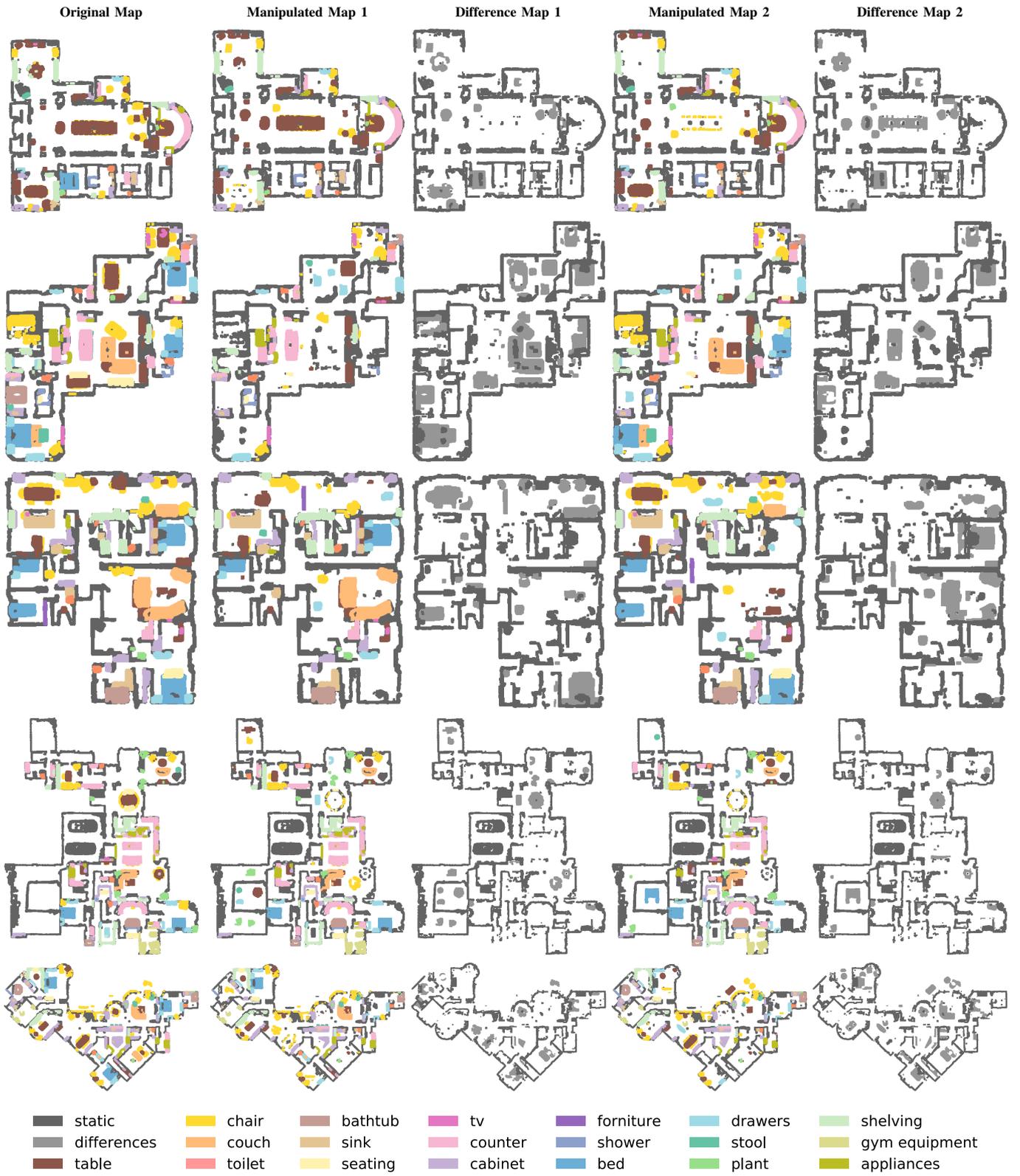


Fig. 9. Generation of alternative states of an environment: original and sample manipulated semantic maps with relative difference maps.