

CLIP on Wheels: Zero-Shot Object Navigation as Object Localization and Exploration

Samir Yitzhak Gadre[◊] Mitchell Wortsman[†] Gabriel Ilharco[†]

Ludwig Schmidt[†] Shuran Song[◊]

Abstract

Households across the world contain arbitrary objects: from mate gourds and coffee mugs to sitars and guitars. Considering this diversity, robot perception must handle a large variety of semantic objects without additional fine-tuning to be broadly applicable in homes. Recently, zero-shot models have demonstrated impressive performance in image classification of arbitrary objects (i.e., classifying images at inference with categories not explicitly seen during training). In this paper, we translate the success of zero-shot vision models (e.g., CLIP) to the popular embodied AI task of object navigation. In our setting, an agent must find an arbitrary goal object, specified via text, in unseen environments coming from different datasets. Our key insight is to modularize the task into zero-shot object localization and exploration. Employing this philosophy, we design CLIP on Wheels (CoW) baselines for the task and evaluate each zero-shot model in both Habitat and RoboTHOR simulators. We find that a straightforward CoW, with CLIP-based object localization plus classical exploration, and *no additional training*, often outperforms learnable approaches in terms of success, efficiency, and robustness to dataset distribution shift. This CoW achieves 6.3% SPL in Habitat and 10.0% SPL in RoboTHOR, when tested zero-shot on *all* categories. On a subset of four RoboTHOR categories considered in prior work, the same CoW shows a 16.1 percentage point improvement in SUCCESS over the learnable state-of-the-art baseline.

1 Introduction

The set of objects humans care about is inherently broad. While object-centric image benchmarks like ImageNet-1k [57], ImageNet-21k [20], MS-COCO [40], and LVIS [28] contain considerable diversity in the object categories they represent, they cannot capture all objects that people might care about in daily life. For example, a first-aid kit can suddenly become relevant in the event of a medical emergency, while a sparkly pacifier might be useful to soothe a crying baby. Hence, we would like household robots to identify and find *arbitrary objects* conditioned on human needs. Furthermore, such object-navigational capabilities should be enabled out-of-the-box, even in *unseen environments* and *without additional training*.

While object navigation settings are well studied in the literature (e.g., [78, 69]), it is typical to fix object categories before training. Existing testing protocols then consider only seen categories in unseen rooms. This setup is akin to standard ImageNet-1k, where categories are likewise fixed. A more challenging, but also more applicable, version of this task is *zero-shot object navigation*, where both target categories and environments are unseen by the agent. This setting naturally takes into account the rich diversity of objects and scenes present in everyday life.

A promising direction for zero-shot object navigation is to use vision models pre-trained on large, heterogeneous datasets as a foundation. For example, CLIP models [53] are pre-trained by matching images from the internet with their textual captions. After pre-training, these models can be used for zero-shot image classification, using human specified captions for the target classes (e.g., “an image of a plant”). Given a

[◊]Columbia University, [†]University of Washington. Correspondence to sy@cs.columbia.edu.

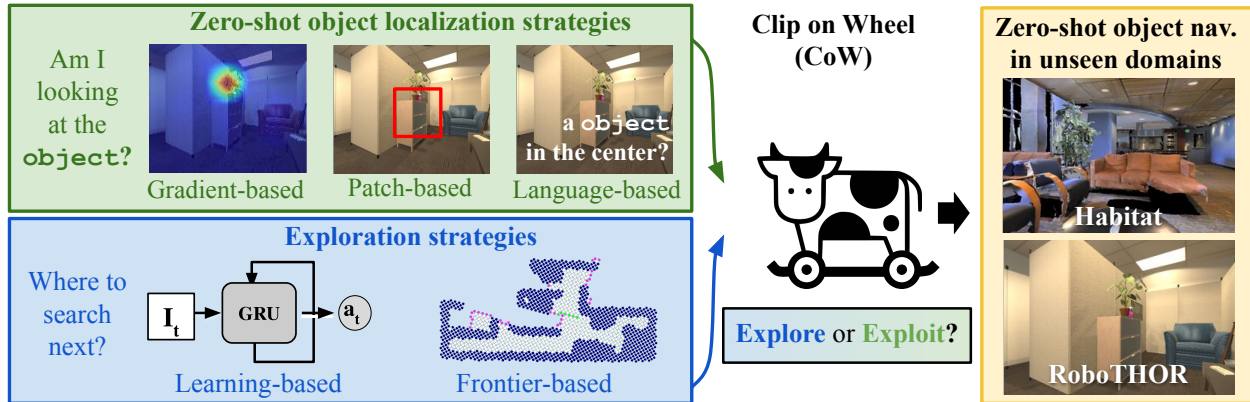


Figure 1: **Overview.** We present the CLIP on Wheels (CoW) framework, which uses an object localization strategy to determine if an object goal is in view and a policy to explore when the goal is not detected. We experiment with three different object localization strategies and two different exploration strategies. Ultimately, our CoWs roam in unseen domains (Habitat and RoboTHOR) to find unseen objects.

test image, the model computes similarity scores between the image and each caption, selecting the caption with the highest score to determine the class label. CLIP-like training has been shown to scale favorably as data and model size increase, becoming increasingly successful for downstream zero-shot image classification [53, 33, 51]. The resulting models have also shown unprecedented robustness to data distribution shifts when compared to standard models trained on a target distribution [53, 71].

However, directly adapting the CLIP model to our task poses several challenges: (1) CLIP does not directly provide a detailed spatial localization of the object when it is in view. This information is critical to inform an agent’s low-level actions. (2) CLIP only processes static images, lacking memory or other mechanisms to guide the agent’s exploration in the environment. (3) Standard fine-tuning of CLIP models can significantly reduce their robustness to data distribution shifts and generality [53, 2, 71, 70]—key strengths of CLIP we want to leverage.

To address these limitations, our *CLIP on Wheels (CoW)* framework factorizes the zero-shot object navigation task into zero-shot object localization and exploration. To preserve CLIP’s generalization capabilities, our framework considers three designs to directly utilize the CLIP model as an object localizer without any fine-tuning. In addition, we also investigate two exploration strategies, considering both learnable and traditional methods. Our algorithm for combining object localization and exploration is similar to one a human may employ: when a CoW has not located the target object, it explores, otherwise, it moves to the target. Our proposed framework is summarized in Fig. 1.

To validate our framework, we study *many* CoWs. Since our method is zero-shot, we evaluate CoWs on both the Habitat [5] and RoboTHOR [17] object navigation domains. We directly test on *all* object categories provided in each benchmark. This is a departure from prior work that usually considers only one simulator, trains different models for each simulator, or considers only a subset of dataset objects at zero-shot test time. We find empirically that CLIP [53] with Grad-CAM [61] object localization and a modified frontier based exploration technique [73]—*with no additional training or fine-tuning*—is an effective zero-shot object navigator. This breed of CoW outperforms competing approaches in terms of (1) success on zero-shot object navigation, (2) efficiency in solving the task, and (3) robustness to data distribution shift (i.e., Habitat v.s. RoboTHOR). This CoW achieves 6.3% SPL in Habitat and 10.0% SPL in RoboTHOR. We also compare against an end-to-end learnable prior work, which tests zero-shot on a subset of four RoboTHOR categories. On this split, our CoW outperforms the state-of-the-art baseline by 16.1 percentage points in navigation SUCCESS.

2 Related Work

Mapping and exploration. Effectively exploring a 3D environment with a mobile robot is a long-standing problem in vision and robotics. Classical methods often decompose the task into map reconstruction [31, 46, 47, 29, 62], agent localization [16, 48, 18], and planning [37, 68], developing mature algorithms for each of these components. While these frameworks are used widely for different robotics applications [21], they are known to be fragile in the presence of uncertainty [41]. Recent work investigates learnable alternatives for exploration [50, 7, 55, 49]. In this line of work, the agent is typically trained end-to-end with self-supervised rewards (e.g., curiosity [50]). Some researchers also propose supervised rewards to encourage exploration [65, 24]. These learning-based methods typically require less manual engineering; however, this comes at the cost of millions of training steps and tedious reward engineering. Since learning-based methods usually optimize performance on a specific distribution, they may be susceptible to distribution shifts and their performance may degrade substantially when deployed in domains that differ from the training distribution. We test both classical and learning based exploration strategies in the context of zero-shot object navigation and study their advantages and disadvantages in our framework.

Goal-conditioned navigation. Apart from open-ended exploration, many navigation tasks are “goal-conditioned,” where the agent needs to navigate to a given position (i.e., point goal [58, 72, 74, 67, 11, 30]), a view of the environment (i.e., image goal [78, 43, 56]), or a specified object category (i.e., object goal [69, 10, 9, 75, 39, 1, 64]). Our task falls into the category of object goal navigation. However, different from the traditional definition, our task does not assume a fixed set of object categories and requires the algorithm to handle unseen object categories at test time without any additional training (i.e., zero-shot inference).

Navigation robustness to domain shift. In different navigation benchmarks (e.g., Habitat [59] and RoboTHOR [17]), an agent could experience drastically different object and layout distributions, visual appearance and style, as well as different noise patterns from sensors and actuators. While prior work studies robust navigation by considering various sensor and actuation noise [12], the testing environments are limited to the RoboTHOR simulator. Khandelwal *et al.* [36] train separate models for Habitat and RoboTHOR to handle the domain shifts between datasets and further limit their zero-shot testing to four RoboTHOR categories (with training conducted on the other eight RoboTHOR categories). In contrast, we evaluate *single models* trained on zero object categories and test on all 12 RoboTHOR categories, additionally testing on all 21 Habitat categories. Most similar to our work is that of Gordon *et al.* [25] who decouple navigation into perception and acting, enabling faster transfer to new simulated environments. However, their model is still updated in new environments and is not capable of navigating to unseen objects.

Zero-shot inference. Zero-shot techniques have become increasingly popular in image classification [53, 33, 51], object detection [4, 54, 19, 38, 77, 42, 26], and semantic segmentation [6, 35, 32, 3, 14, 76]. Recent zero-shot models for image classification include CLIP [53], ALIGN [33], and BASIC [51]. These methods pre-train on image-text pairs from the internet, where the objective is to match an image with its corresponding caption. These models can then be used for zero-shot image classification by constructing candidate captions for each of the target classes directly at test time. For instance, a zero-shot classifier which discriminates between cows and plants may be instantiated using the captions “a photo of a cow” and “a photo of a plant”. The predicted class will be the caption that is better aligned with the test image. We study how to adapt such models to enable zero-shot object navigation with no additional training.

3 Approach

This section defines the zero-shot object navigation task along with the environment and embodiment details we consider (Sec. 3.1). Next, we provide an overview of our approach (Sec. 3.2). In the following sections, we describe various object localization (Sec. 3.3) and exploration (Sec. 3.4) strategies. We then present simple recipes to combine these ingredients to build CLIP on Wheels models (CoWs), which are zero-shot object navigators (Sec. 3.5).

3.1 Zero-Shot Object Navigation

Task details. We define the task of navigating to unseen goal objects in unseen domains, modifying the notation from [69]. Let \mathcal{O} denote a set of unseen object categories and instances (e.g., `plant`, `blue cow stuffed animal`, etc.) that an agent should find. Let \mathcal{S} denote the set of unseen scenes, which an agent should navigate. Let p_0 describe the initial pose of the agent in a scene. A navigation episode $\tau \in \mathcal{T}$ can then be written as a tuple $\tau = (s, o, p_0)$, where $s \in \mathcal{S}, o \in \mathcal{O}$. Starting at p_0 , an embodied agent is initialized to find target o . The agent has some action space \mathcal{A} and receives observation and sensor readings I_t . At each timestep t , the agent executes an action $a \in \mathcal{A}$, issuing a special action `STOP` to terminate the navigation episode. If the agent is within c units of the object o and o is visible, the episode is considered a success. We stress that neither \mathcal{O} nor \mathcal{S} is explicitly modeled during training and hence \mathcal{T} is a set of zero-shot test tasks.

Environment and embodiment details. We adopt the Habitat-challenge [5] and RoboTHOR-challenge [17] object navigation validation sets as our test sets. The ground-truth for the challenge test sets is not publicly available, motivating this decision. These two domains differ substantially in terms of their appearance. Habitat uses Matterport3D [8] reconstructed scenes while RoboTHOR uses artist generated scenes. Furthermore, the underlying physics engines are also different. For the rest of this section, we present dataset information for Habitat and RoboTHOR respectively. There are a total of 31 object categories ($|\mathcal{O}| = 31$), split 21 / 12. Two categories (`plant` and `TV`) overlap between the domains. There are 26 test time scenes ($|\mathcal{S}| = 26$), split 11 / 15 and 3,995 episodes ($|\mathcal{T}| = 3,995$), split 2,195 / 1,800. Agents receive RGB-D egocentric views. All agents have discrete actions $\{\text{MOVEFORWARD}, \text{ROTATERIGHT}, \text{ROTATELEFT}, \text{STOP}\}$ ($|\mathcal{A}| = 4$). The move action advances the agent by 0.25m, while rotation actions pivot the agent’s camera by 30° . The agent is a LoCoBot [27]. The goal distance threshold is $c = 1.0\text{m}$. Both domains are setup with noise that is faithful to the original CVPR 2021 object navigation challenges. For Habitat this means no actuation noise, but considerable noise in depth and visual images as Matterport3D reconstructions contain many artifacts. For RoboTHOR, this means gaussian actuation noise, but no depth sensor noise.

3.2 CLIP on Wheels Overview

To tackle the challenging zero-shot object navigation setting, we explore a modular approach that (1) adapts zero-shot classification models like CLIP to make them suitable for object localization and (2) employs exploration strategies to navigate a space, thereby generating many egocentric views.

We call our approach *CLIP on Wheels (CoW)*. We make use of zero-shot pre-trained models instead of fitting models to Habitat or RoboTHOR directly, which are more limited in visual diversity than the internet image distribution CLIP was trained on. In effect, using the pre-trained model without additional training allows us to keep it zero-shot, while exploration allows us to bring the model into navigation settings. A CoW, then, should explore while it is not confident about the target object localization. When it is confident, it should take actions to bring it towards the target, avoiding obstacles along the way.

3.3 Object Localization Modules

Successful object navigation requires localizing objects and making plans to move towards them. Since our CoWs are enabled with depth sensors, back-projected regions of the image corresponding to high object saliency give natural targets for navigation planning. To incorporate semantic localization, we consider three strategies: gradient-based, patch-based, and language-based. Qualitative saliency maps of each strategy are visualized in Fig. 2 for the `plant` category in both Habitat and RoboTHOR. An example of back-projecting saliency into an estimated map can be seen in Fig. 3 (left). While sophisticated methods for zero-shot localization (e.g., bounding box detection, semantic segmentation) exist, they often require additional learned components (e.g., [26, 76]). In the spirit of using CLIP zero-shot, we leave investigation of incorporating such methods to future work.

Review of CLIP for zero-shot classification. Recall that CLIP is contrastively pre-trained on a dataset of 400M image-caption pairs. The model is encouraged to align feature representations of corresponding

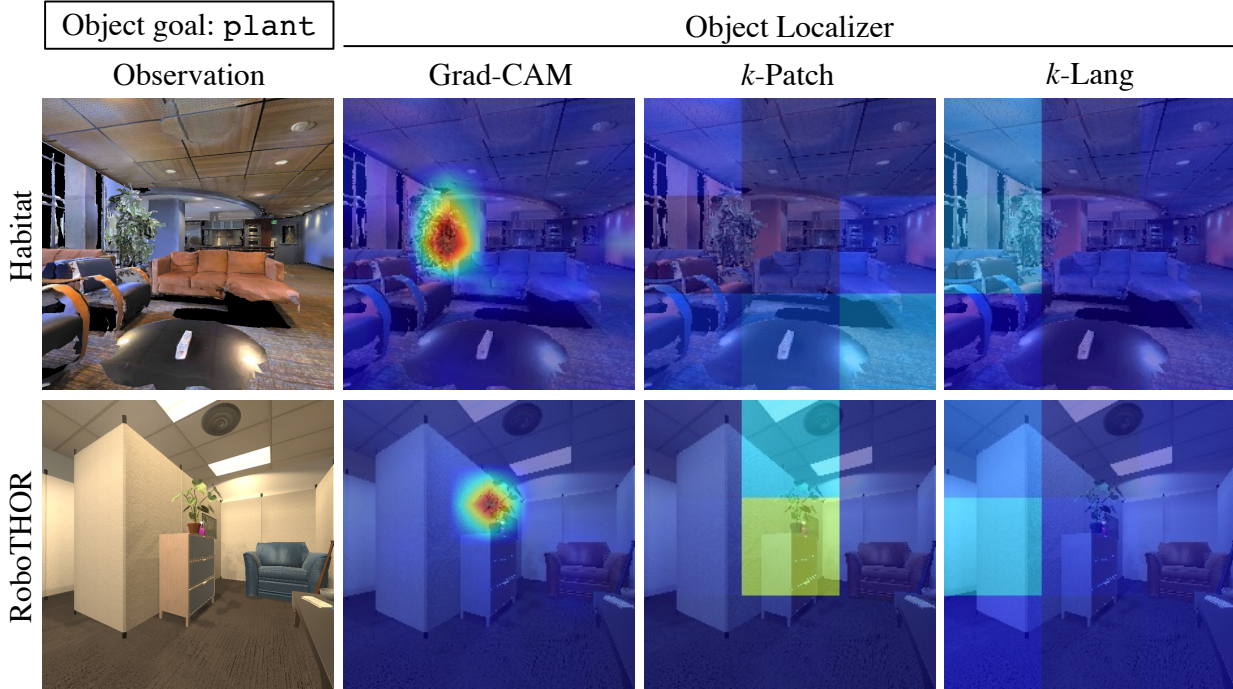


Figure 2: **Zero-Shot Localization.** Visualization of the three localization strategies considered, with CLIP ViT-B/32 as the underlying model. Warmer colors indicate higher object saliency. Here, Grad-CAM [61] qualitatively gives the best localization of the `plant` in both environments.

images and captions, while pushing these representations away from non-matching images and captions. At test time, a classifier is created by procedurally generating captions for each target class as “mock” captions. Given an image, CLIP extracts an image feature, and classifies it with the caption corresponding to the most similar caption feature. While incredibly powerful, CLIP does not directly give a way to localize *where* in the image the target is located, thus motivating our investigation.

Grad-CAM [61, 13]. We first consider Gradient-weighted Class Activation Mapping (Grad-CAM) [61] for object localization in RGB images. Using the target text feature (i.e., a CLIP caption feature representing the target object category), this technique uses the gradients of the CLIP image encoder with respect to the target to create a saliency map over the image. Empirically, Grad-CAM highlights areas that are deemed important for the image feature to be aligned with its text label. For transformer-based architectures (i.e., ViT [22]), we adopt the modifications developed by Chefer *et al.* [13]. Fig. 2 (left) shows example saliency maps produced by Grad-CAM. We also observe that when targets are not in view, Grad-CAM give low saliency predictions, which is useful for preventing false positive detections. See Appx. A for a more formal treatment of Grad-CAM and additional visualization. To turn the saliency map into a discrete prediction of the region of interest, we first scale maps with a multiplicative constant, clip between zero and one, and threshold the predictions. For more details on hyperparameters see Appx. B.

k -Patch. In contrast to Grad-CAM, this strategy only requires forward passes of the CLIP model, without any backward gradient information. The central idea is to discretize the image into k smaller patches, and run inference on each of them. For instance, if the goal is to locate a “`plant`” in an image, we can divide the image into 3×3 smaller patches, extract image features for each patch, and compute similarity with the CLIP language feature representation of the caption “`a photo of a plant.`” The similarity values can be converted to a distribution over patches using standard `SoftMax` to get a probability that each patch is in fact a “`plant`”. Based on the observation that CLIP models are well calibrated [44] (i.e., CLIP confidence is indicative of CLIP predictive accuracy), a high patch probability should correlate with the object being

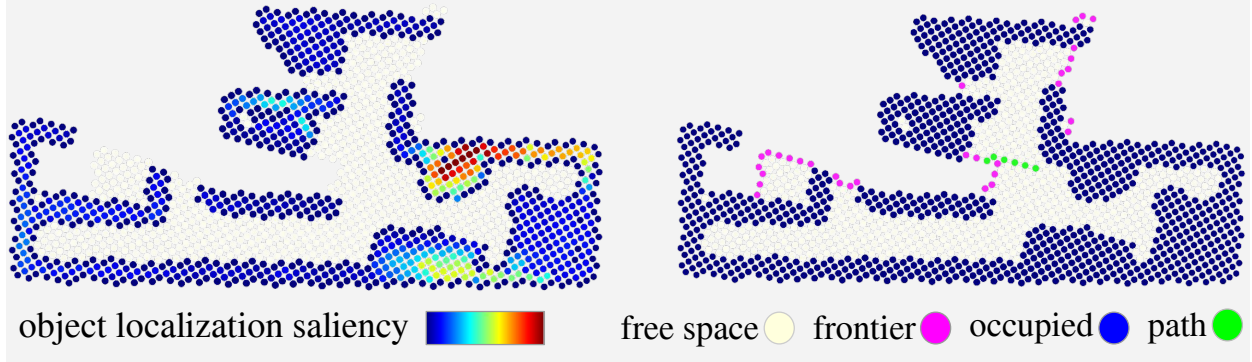


Figure 3: **Mapping.** Top-down map created from egocentric depth observations as the agent explores the space. **(left)** Back-projected Grad-CAM [61] object localization saliency using depth and agent pose estimation. **(right)** A visualization of Frontier Based Exploration (FBE) [73] showing the planned path to the next frontier.

in that patch. See Appx. B for the confidence threshold. In practice, patches can be batched, allowing for faster inference. An illustration of this approach is provided in Fig. 2 (middle) highlighting a potential pitfall when an object is split between patches.

***k*-Lang.** Similarly to *k*-Patch, *k*-Lang does not require any modifications to CLIP. However, instead of discretizing the image into smaller patches, *k*-Lang uses a single pass over the entire image, but attempts to match it with *k* different captions that capture positional information about the object. Taking our running example of locating a “plant”, the captions could be “a photo of a plant in the top left.”, “a photo of a plant in the bottom right.”, etc. CLIP feature similarity is computed between the image feature and each language feature. Again, SoftMax is computed over the similarity scores and confident predictions for a locale gives a way to recover a patch. See Appx. B for more details. Fig. 2 (right) illustrates this strategy, which appears to miss the “plant” entirely.

3.4 Exploration Modules

Successful object localization identifies the target when it is in view. However, in many egocentric images, the target may *not* be visible. Hence, a successful navigator should explore a space to generate views for the object localization algorithm. Towards this goal, we first discuss depth based mapping. We then discuss frontier-based and learning-based methods for action selection. Finally, we discuss a simple strategy for action failure checking, which can be incorporated into both exploration strategies.

Depth based mapping. As a CoW moves around a scene, we would like to construct a top-down voxel map on-the-fly as shown in Fig. 3. Such a map should capture free, occupied, and unknown space, which provides valuable information about where an agent can move. To create this map, we first estimate the pose of the CoW’s coordinate frame C with respect to a reference frame W .

Consider a CoW at initial position and orientation p_0 . The initial pose of the CoW in relation to frame W is then ${}^W T_{C_0} = I$, where I is the identity SE(3) transform. Since a CoW knows the intended consequences of its actions (e.g., MOVEFORWARD should result in a translation of 0.25m), each action can be represented as a fixed pose delta transform ${}^{C_t} \hat{T}_{C_{t+1}}$, which models an action transition. To estimate the CoW pose at timestep $t + 1$, we compute ${}^W \hat{T}_{C_{t+1}} = {}^W \hat{T}_{C_t} \cdot {}^{C_t} \hat{T}_{C_{t+1}}$. Due to sensor-noise and potential for failure of executed actions, this estimate may not be accurate. Furthermore, since transforms are simply concatenated, this formulation faces the problem of pose drift. We discuss strategies for managing such complexity later in this section, but use this formulation for clarity.

We utilize the current estimated pose ${}^W T_{C_t}$, camera intrinsic matrix K , current depth observation D_t , standard back back-projection of D_t , and pose concatenation to register new observations in a global 3D map (in coordinate frame W). Since navigation is mostly concerned with free v.s. occupied space, we do not keep the whole 3D map. Rather we project the map on the ground-plane using the known agent height and

assumed down gravity direction. Points already near the ground-plane are considered free space, while other points are considered occupied as shown in Fig. 3. Additionally, we voxelize the map to keep it parsimonious. Voxel resolution is set to 0.125m in our case, which is half the CoW MOVEFORWARD distance of 0.25m.

Using the saliency map computed via methods discussed in Sec. 3.3, it is also possible to register predictions into the voxel map (i.e., each 3D point before vocalization also has an object saliency value associated with it). This allows for a way to “lift” the 2D CLIP-based object localization to 3D. Because the map is projected to the ground-plane, we keep track of only the max saliency value for each voxel. See Fig. 3 (left) for visualization. Once a target region of interest is registered in the map, it becomes trivial to plan a path towards that area of the map.

Frontier based exploration (FBE) [73]. The top-down map can be extended to implement a simple exploration method introduced in 1997 by Yamauchi [73]. At initialization the agent spins in a 360° fashion to comprehensively visualize its surroundings and initialize the map. Consider the boundaries between free space and unknown space (visualized in purple in Fig. 3 (right)). Intuitively, by moving towards these regions, called *frontiers*, a navigator should discover new pieces of the map. With the voxel grid of free space, the agent computes the shortest path (e.g., breath first or A* path) to the nearest frontier (shown in green in Fig. 3 (right)). Because a CoW’s action space is not grid aligned, a low-level planner additionally searches for sequences of actions that make progress towards way-points along the shortest voxel path. The low-level planner could also be implemented with a continuous action space, highlighting the extensibility of this method. Once the navigator reaches a frontier, it moves to the next closest frontier.

The map is updated at every timestep, so noisy pose estimates contribute to inaccuracies in the map, which also makes exploration planning harder. For example, narrow passageways may “collapse,” with opposite walls fusing together in the voxel grid. A CoW may then think these hallways are not traversable. To circumvent such problems we simply reinitialize the map when the CoW does not believe it can reach any frontier (i.e., no paths through free space exist between the current CoW voxel and frontiers). Incorporating more sophisticated recovery strategies is left to future work.

Learnable exploration. In addition to frontier-based exploration, we investigate learnable alternatives. We consider an architecture and reward structure similar to other work in embodied AI (e.g., [65, 36]). Specifically we adopt a frozen, pre-trained vision backbone outfit with a trainable GRU [15] and linear actor, critic heads. Inspired by the recent success of frozen CLIP backbones for embodied AI [36], we adopt a frozen CLIP ViT-B/32 backbone as our vision encoder. Hence, our network takes RGB input at 224×224 resolution. We train agents independently in Habitat [59] and RoboTHOR [17] simulation environments for 60M steps each, using DD-PPO [60, 67] and the AllenAct [66] framework. For hyperparameters please see Appx. C. We employ a simple count-based reward [63], which self-supervised exploration methods often attempt to approximate (e.g., [7]). During training, an agent receives a reward of 0.1 for visiting a previously unvisited voxel location (at 0.125m resolution), and a step penalty of -0.01. We dub the resulting networks the Habitat-Learnable Explorer (Hab.-Learn) and RoboTHOR-Learnable Explorer (Robo.-Learn), corresponding to the domains in which they were trained. The training datasets are disjoint from the downstream zero-shot navigation test sets. For more detail on training datasets see Appx. D.

Action failure checking. For both learnable and heuristic explorers, actions may fail. For map based explorers, an obstacle might be below the field of view and hence not captured as occupied space. For learnable explorers, the output policy may heavily favor a failed action (e.g., MOVEAHEAD), resulting in a CoW getting stuck trying to roll into occupied space. To improve the action success of our CoWs (empirically verified in Sec. 4.2), we employ a simple strategy. We compute the absolute difference between depth channels in the observations $\Delta_{i,i+1} = |D_i - D_{i+1}|$. We then compute the mean $\mu(\Delta_{i,i+1})$ and standard deviation $\sigma(\Delta_{i,i+1})$ as representative statistics. These quantities have interpretable meaning in meters. Since actions should move the agent forward by approximately a fixed distance or rotation, we can then set reasonable thresholds for μ , σ below which we can be confident actions failed. See Appx. E for threshold values. In modern robot navigation systems, on-board pose estimation and bumper sensors for obstacle avoidance are nearly ubiquitous. In the real world it is possible to implement action failure checking with even more sensor information.

3.5 A CoW Emerges from Object Localization and Exploration

As show in Fig. 1, to build a CoW, we first choose one object localizer: Grad-CAM, k -Patch, or k -Lang. All CoWs have access to depth sensors and hence can construct 3D object saliency maps as shown in Fig. 3 (left). A CoW also gets an explorer: Frontier-based exploration (FBE), Habitat-learnable exploration (Hab.-Learn), or RoboTHOR-learnable exploration (Robo.-Learn). When a CoW has not localized an object it uses its exploration module to roam. Once the localizer fires for the target object (i.e., \max saliency exceeds a threshold), the CoW switches to exploit mode and navigates directly to the region of interest by planning within the voxel map. When it reaches the target it issues STOP.

4 Experiments

In this section we start by evaluating our zero-shot object localizers (Sec. 4.1) and explorers (Sec. 4.2) independently on Habitat and RoboTHOR environments to build intuition about their performance. For modules that generate reasonable performance, we construct CLIP on Wheels models (CoWs) that are evaluated on the end task of zero-shot object navigation (Sec. 4.3).

4.1 Object Localization

We evaluate three zero-shot localization strategies (Grad-CAM, k -Patch, and k -Lang presented in Sec. 3.3) on static images from Habitat and RoboTHOR to quantify the degree to which CLIP can localize objects.

Datasets. For the 21 object categories in Habitat and 12 object categories in RoboTHOR, we render datasets, with 1,279 and 540 images respectively. The Habitat images contain an average of 2.98 object categories per image while the average in RoboTHOR is 2.34. Each image contains at least one object category. For more details on the image sampling procedure please see Appx. F.

Algorithms. We consider Grad-CAM, k -Patch, and k -Lang approaches presented in Sec. 3.3. We consider $k = 9$, resulting in 9-Patch and 9-Lang strategies visualized in Fig. 2. To turn saliency maps into localization regions of interest, we threshold predictions. Hence, each localizer saliency map is ultimately converted to a binary mask prediction. For more details on thresholds, please see Appx. B.

We also ablate two prompt strategies. The first (VG) leverages the domain knowledge that Habitat and RoboTHOR are simulators and may have the appearance of video games. This motivates the prompt: “a photo of a <object> in a video game.” The second (OpenAI) is to use the prompt ensemble of 80 prompts used by Radford *et al.* [53] for ImageNet-1k zero-shot evaluation. When using more than one prompt per category, the idea is to average text representations over all prompt templates for each category. This is empirically shown to improve results [53] on ImageNet-1k. This set of prompts includes potentially helpful templates like “a photo of a hard to see <object>.” However, it also includes templates like “a tattoo of the <object>.” that do not intuitively describe Habitat and RoboTHOR data.

Metrics. We evaluate zero-shot object location performance using $F1$ -score. Recall that object saliency maps like those in Fig. 2 are converted to binary localization mask predictions by thresholding on a constant. For each image, there are some object categories \mathcal{O}^+ that appear in that image and other categories \mathcal{O}^- that do not. We consider a predicted binary localization mask $M_{\text{pred}}^{o^+}$ a true positive (TP) if $\text{score}^+ = \sum(M_{\text{pred}}^{o^+} \odot M_{\text{GT}}^{o^+}) / \sum M_{\text{pred}}^{o^+} > 0.5$, where $o^+ \in \mathcal{O}^+$, $M_{\text{GT}}^{o^+}$ is the ground truth mask, the summation is over all binary pixel values, and \odot is an element-wise product that gives binary mask intersection. This is a more lenient measure than traditional jaccard index; however, also more applicable for our navigation setting where only part of the object needs to be identified correctly to provide a suitable navigation target. False positives (FP) arise when $0 < \text{score}^+ \leq 0.5$ or $\sum M_{\text{pred}}^{o^-} > 0$, where $o^- \in \mathcal{O}^-$. False negatives (FN) arise when $\sum M^{o^+} = 0$. $F1$ is computed in the standard way: $\frac{\text{TP}}{\text{TP}+0.5(\text{FP}+\text{FN})}$. In each domain $F1$ is computed per category and then averaged over the categories (i.e., macro $F1$). This yields $F1^H$ for Habitat and $F1^R$ for RoboTHOR.

Localizer	Grad-CAM		9-Patch		9-Lang	
Prompt set	VG	OpenAI	VG	OpenAI	VG	OpenAI
$F1^H$	0.074	0.084	0.029	0.036	0.000	0.000
$F1^R$	0.093	0.103	0.018	0.021	0.000	0.000

Table 1: **Zero-shot object localization F1 comparisons.** We study the effects of different object localizer and prompt strategies by looking at $F1$ -score ($F1^H$ in Habitat and $F1^R$ in RoboTHOR), which takes into account both precision and recall. 9-Lang achieves poor performance, while both Grad-CAM and 9-Patch appear to work in some cases. Additionally, the 80 prompt ensemble (OpenAI) outperforms the “video game” prompt (VG).

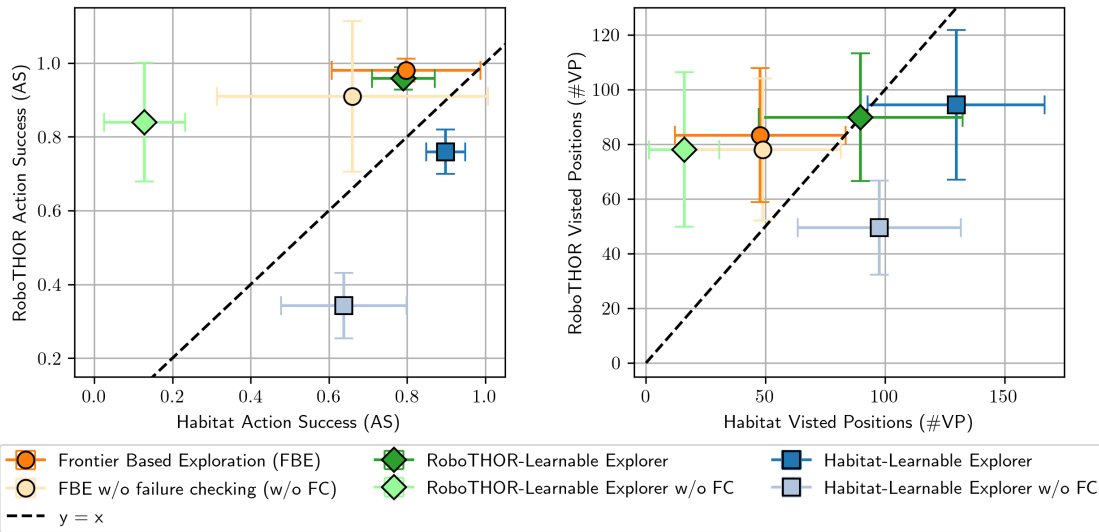


Figure 4: **Exploration Results.** Agents are rolled-out for 250 steps from three random initializations in Habitat and RoboTHOR in all testing rooms. Both learnable methods experience a performance drop out-of-distribution; however, this can be mitigated via failure checking (FC) with the added benefit of in-distribution gains. **(left)** Action Success AS decreases without failure checking (w/o FC). **(right)** Both learned explorers, Habitat-Learnable w/o FC and RoboTHOR-Learnable w/o FC perform comparably or better than Frontier-based exploration (FBE) when considering in-distribution performance.

Effects of different object localizers. Looking Tab.1 (right), we see that 9-Lang achieves poor performance across different benchmarks and is not suitable as a localizer. This result is perhaps expected; intuitively, people seldom localize objects spatially in internet captions (i.e., data distribution CLIP was pre-trained on). Grad-CAM outperforms 9-Patch; however, we consider both localizers downstream (Sec. 4.3).

Effects of text prompts. In Tab. 1, we find that the 80 prompt ensemble templates (OpenAI) outperform the “video game” template (VG); however, we consider both prompt strategies downstream (Sec. 4.3).

4.2 Exploration

In this set of experiments, we compare the performance of frontier-based and learning-based methods for exploration.

Datasets. To evaluate exploration effectiveness, we sample three random starting points in each of the 11 Habitat and 15 RoboTHOR zero-shot testing rooms. CoWs are rolled out for 250 steps for each exploration episode.

Algorithms. We consider two types of exploration strategies: frontier-based and learning-based. This results in three instances of explorers: a frontier based explorer, an explorer learned in Habitat, and an explorer learned in RoboTHOR. For each of the three agents we evaluate without failure checking (w/o FC) and with FC, which is the default setting. This amounts to six exploration models in our test-bed. We say a learned agent is out-of-distribution (OOD) when it is trained in one simulation environment and tested in another. Analogously, an agent is in-distribution (ID) when evaluated in the same environment as it is trained.

Metrics. We use the following metrics to evaluate the effectiveness of different exploration policies. Since metrics are evaluated independently in each simulation domain, we are able to use these metrics to also evaluate differences between ID and OOD performance.

- Number of visited position (#VP): the average number of positions that an agent visits. Counts are recorded over 250 steps. We discretize positions by normalizing by 0.125m (i.e., half of the agent’s step length) and rounding to the nearest integer. Rotation diversity does not explicitly help this metric; however, rotating away from obstacles can implicitly help for future timesteps. This is a common metric in the exploration literature (e.g., [49]), where agents visiting more positions are considered better.
- Action Success (AS): the fraction of actions that were successfully executed averaged over all timesteps and considering all evaluation episodes. This metric gives insight into how likely the agent is to bump into obstacles, which is of critical importance in real world deployment.

How do explorers trained in different domains compare? Perhaps unsurprisingly, we find explorers perform better ID than OOD, as seen in Fig. 4. Interestingly, failure checking (FC) improves both ID and OOD performance, with larger gains OOD.

How do heuristic and learnable explorers compare? In Fig. 4 (left) we see that FBE performs comparably to the Learnable Explorers. In Fig. 4 (right) both Learnable-Explorers visit more positions than FBE, ID and OOD, suggesting they are reasonable points of comparison for the downstream task.

4.3 Zero-Shot Object Navigation

To validate our CoW zero-shot object navigators we evaluate them on both Habitat and RoboTHOR domains, using the official object navigation val. sets as our test sets—the full test sets are not publicly available. All CoWs estimate voxel maps on-the-fly and employ the same target-oriented planning strategy when a localizer fires. This allows for fair-comparison across CoWs.

Metrics. We use the following metrics for evaluating agents. For every metric, higher is better.

- SUCCESS: the fraction of episodes where the agent executes the STOP action within 1.0m of the target object.
- Success weighted by inverse path length (SPL): success is weighted by the oracle shortest path length and normalized by the actual path length [5]. Hence this metric points to the success efficiency of the agent. This is the main metric in the CVPR 2021 object navigation challenges.
- Action Success (AS): the fraction of actions that were successfully executed, same as the exploration metric in Sec. 4.2.

Overall performance. Based on the results in Tab. 2, we see that the CoW composed of FBE, OpenAI prompts, and Grad-CAM performs well on both Habitat and RoboTHOR, achieving the highest SPL (the main challenge metric) in both domains. The second best explorer in Habitat in terms of SPL is Hab.-Learn, which employs an exploration agent trained in Habitat. The Hab.-Learn CoWs marginally beats the FBE

Exploration	CoW breeds		Fine-tuning steps	Habitat			RoboTHOR		
	Prompts	Localizer		SPL↑	SUCCESS↑	AS↑	SPL↑	SUCCESS↑	AS↑
FBE [73]	OpenAI	Grad-CAM [61]	0	0.063	0.111	0.828	0.100	0.163	0.980
FBE [73]	VG	Grad-CAM [61]	0	0.051	0.092	0.832	0.095	0.161	0.978
FBE [73]	OpenAI	9-Patch	0	0.047	0.088	0.834	0.093	0.154	0.981
FBE [73]	VG	9-Patch	0	0.042	0.084	0.826	0.091	0.148	0.981
Hab.-Learn	OpenAI	Grad-CAM [61]	60M	0.044	0.115	0.806	0.060	0.148	0.756
Hab.-Learn	VG	Grad-CAM [61]	60M	0.044	0.113	0.801	0.052	0.142	0.755
Robo.-Learn	OpenAI	Grad-CAM [61]	60M	0.032	0.086	0.787	0.084	0.147	0.949
Robo.-Learn	VG	Grad-CAM [61]	60M	0.029	0.078	0.790	0.080	0.142	0.947

Table 2: **Zero-shot object navigation in Habitat and RoboTHOR.** CoW w/ frontier-based exploration (FBE) + Grad-CAM for CLIP localization performs favorably compared to alternatives that employ learned exploration in Habitat (Hab.-Learn) and RoboTHOR (Robo.-Learn). Fine-tuning steps refers to the number of steps the learnable models were fine-tuned for exploration. The best CoW in terms of SPL is highlighted in blue.

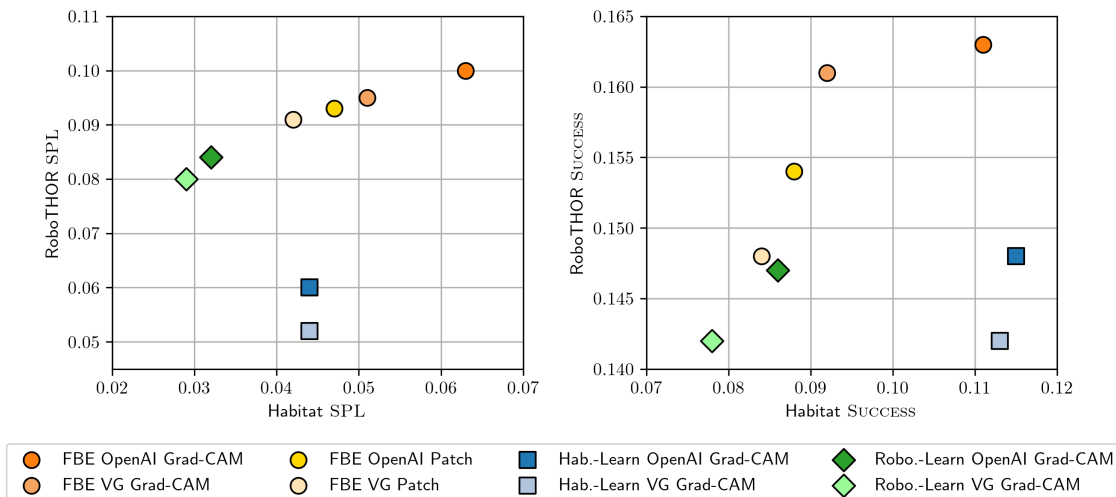


Figure 5: **Scatter view of Tab. 2** All points lie above the $y = x$ line. **(left)** FBE points are most robust to the domain shift. **(right)** Similar trends with SUCCESS as with SPL. One Hab.-Learn CoW outperforms FBE-based CoWs in Habitat on this metric; however, this model is also less robust.

CoWs in terms of SUCCESS in Habitat, suggesting there can be some advantages to training CoWs on target data. For category level results, see Appx. H.

Safe navigation. We again see that FBE enabled CoWs outperform competing methods in Tab. 2, achieving an average of 83% AS in Habitat and 98% AS in RoboTHOR. This outperforms all learned exploration CoWs, which tend to execute successful actions at a lower rate, even with equal access to the depth based failure checking discussed in Sec. 3.4.

Robustness to domain shift. Looking at Tab. 2 we see, across all metrics, learnable agents perform better ID than OOD. FBE CoWs show the best performance in both domains. To make these trends more clear, we provide a scatter-plot view of Tab. 2 in Fig 5.

Exploration	CoW breeds		Fine-tuning steps	RoboTHOR SUCCESS				
	Prompts	Localizer		apple	basketball	plant	television	All
FBE [73]	OpenAI	Grad-CAM [61]	0	0.127	0.160	0.413	0.267	0.242
FBE [73]	VG	Grad-CAM [61]	0	0.153	0.167	0.396	0.247	0.240
FBE [73]	OpenAI	9-Patch	0	0.160	0.167	0.240	0.213	0.195
FBE [73]	VG	9-Patch	0	0.147	0.167	0.260	0.153	0.182
Hab-Learn	OpenAI	Grad-CAM [61]	60M	0.173	0.113	0.253	0.267	0.202
Hab-Learn	VG	Grad-CAM [61]	60M	0.167	0.113	0.267	0.200	0.187
Robo-Learn	OpenAI	Grad-CAM [61]	60M	0.147	0.107	0.320	0.213	0.197
Robo-Learn	VG	Grad-CAM [61]	60M	0.127	0.140	0.333	0.193	0.198
EmbCLIP zero-shot [36]			60M	0.147	0.067	0.053	0.060	0.081

Table 3: **Comparison to EmbCLIP zero-shot [36].** Prior work trains on eight RoboTHOR categories and tests on four unseen categories. All CoWs outperform the EmbCLIP zero-shot baseline (see col. “All”). The best CoW, which requires zero steps of navigation training, outperforms the baseline, which was trained for 60M, steps by 16.1pp. The main rows for comparison are highlighted in blue .

Comparison to an end-to-end learnable baseline. To further contextualize our results, we compare our CoW models to prior zero-shot object navigation work: the EmbCLIP zero-shot model [36]. This model uses a frozen CLIP backbone, learnable actor, critic heads, and is trained for 60M steps—similar to our exploration agents. EmbCLIP zero-shot is trained on object navigation episodes for eight RoboTHOR categories, in the same training rooms as our exploration agents, which is disjoint from the testing rooms. The model then tested on four held-out categories. On the four category split, we find that all CoWs outperform EmbCLIP zero-shot by a large margin. The FBE, OpenAI, Grad-CAM CoW achieves an impressive 24.2% success, which is a 16.1pp gain over the state-of-the-art.

5 Conclusion and Future Work

This paper introduces the CLIP on Wheels (CoW) framework, translating the successes of zero-shot image classification models to the popular embodied AI task of object navigation. Without any training, our best CoW model attains 6.3% and 10.0% success weighted by inverse path length (SPL) on zero-shot object navigation in Habitat and RoboTHOR respectively, outperforming analogous baselines that have access to domain-specific data and optimization. Our experiments demonstrate that by modularizing the task into zero-shot object localization and exploration, and selecting the right method for each module, the proposed CoW framework can gain 16.1pp over an end-to-end learning-based method in terms of task success. While our evaluation of single models on both Habitat and RoboTHOR is a step towards evaluating zero-shot object navigators and their robustness to domain shifts, ultimately performance in the real-world matters most. Considering real-world deployment with user specified goals is an exciting directions. Future work should also assess if recent algorithmic progress in Habitat and RoboTHOR actually translates to real-world gains (in the spirit of prior work that indicates it does [34]). Finally, inspired by the findings of Pratt *et al.* [52], future work should consider how modifying robot embodiment (i.e., robot morphology, action space, camera resolution, etc.) can improve navigation performance.

Acknowledgement

We would like to thank Jessie Chapman, Huy Ha, Sachit Menon, and Sarah Pratt for valuable discussions and feedback. This work was supported in part by NSF CMMI-2037101, NSF IIS-2132519, and an Amazon Research Award. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of the sponsors.

References

- [1] Al-Halah, Z., Ramakrishnan, S.K., Grauman, K.: Zero experience required: Plug & play modular transfer learning for semantic visual navigation. arXiv (2022)
- [2] Andreassen, A., Bahri, Y., Neyshabur, B., Roelofs, R.: The evolution of out-of-distribution robustness throughout fine-tuning. arXiv (2021)
- [3] Baek, D., Oh, Y., Ham, B.: Exploiting a joint embedding space for generalized zero-shot semantic segmentation. ICCV (2021)
- [4] Bansal, A., Sikka, K., Sharma, G., Chellappa, R., Divakaran, A.: Zero-shot object detection. ECCV (2018)
- [5] Batra, D., Gokaslan, A., Kembhavi, A., Maksymets, O., Mottaghi, R., Savva, M., Toshev, A., Wijmans, E.: ObjectNav Revisited: On Evaluation of Embodied Agents Navigating to Objects. arXiv (2020)
- [6] Bucher, M., VU, T.H., Cord, M., Pérez, P.: Zero-shot semantic segmentation. NeurIPS (2019)
- [7] Burda, Y., Edwards, H., Storkey, A.J., Klimov, O.: Exploration by random network distillation. ICLR (2018)
- [8] Chang, A.X., Dai, A., Funkhouser, T.A., Halber, M., Nießner, M., Savva, M., Song, S., Zeng, A., Zhang, Y.: Matterport3d: Learning from rgb-d data in indoor environments. 3DV (2017)
- [9] Chang, M., Gupta, A., Gupta, S.: Semantic visual navigation by watching youtube videos. NeurIPS (2020)
- [10] Chaplot, D.S., Gandhi, D., Gupta, A., Salakhutdinov, R.: Object goal navigation using goal-oriented semantic exploration. NeurIPS (2020)
- [11] Chaplot, D.S., Gandhi, D., Gupta, S., Gupta, A.K., Salakhutdinov, R.: Learning to explore using active neural slam. ICLR (2020)
- [12] Chattopadhyay, P., Hoffman, J., Mottaghi, R., Kembhavi, A.: Robustnav: Towards benchmarking robustness in embodied navigation. ICCV (2021)
- [13] Chefer, H., Gur, S., Wolf, L.: Transformer interpretability beyond attention visualization. CVPR (2021)
- [14] Cheng, J., Nandi, S., Natarajan, P., Abd-Almageed, W.: Sign: Spatial-information incorporated generative network for generalized zero-shot semantic segmentation. ICCV (2021)
- [15] Cho, K., van Merriënboer, B., Bahdanau, D., Bengio, Y.: On the properties of neural machine translation: Encoder–decoder approaches. In: SSST@EMNLP (2014)
- [16] Davison, A.J., Murray, D.W.: Mobile robot localisation using active vision. ECCV (1998)
- [17] Deitke, M., Han, W., Herrasti, A., Kembhavi, A., Kolve, E., Mottaghi, R., Salvador, J., Schwenk, D., VanderBilt, E., Wallingford, M., Weihs, L., Yatskar, M., Farhadi, A.: Robothor: An open simulation-to-real embodied ai platform. CVPR (2020)
- [18] Dellaert, F., Fox, D., Burgard, W., Thrun, S.: Monte carlo localization for mobile robots. ICRA (1999)
- [19] Demirel, B., Cinbis, R.G., Ikizler-Cinbis, N.: Zero-shot object detection by hybrid region embedding. BMVC (2018)
- [20] Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. CVPR (2009)

- [21] DeSouza, G.N., Kak, A.C.: Vision for mobile robot navigation: A survey. *IEEE transactions on pattern analysis and machine intelligence* (2002)
- [22] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR* (2020)
- [23] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv* (2020)
- [24] Gadre, S., Ehsani, K., Song, S., Mottaghi, R.: Continuous scene representations for embodied ai. *CVPR* (2022)
- [25] Gordon, D., Kadian, A., Parikh, D., Hoffman, J., Batra, D.: Splitnet: Sim2sim and task2task transfer for embodied visual navigation. *CVPR* (2019)
- [26] Gu, X., Lin, T.Y., Kuo, W., Cui, Y.: Zero-shot detection via vision and language knowledge distillation. *arXiv* (2021)
- [27] Gupta, A.K., Murali, A., Gandhi, D., Pinto, L.: Robot learning in homes: Improving generalization and reducing dataset bias. *NeurIPS* (2018)
- [28] Gupta, A., Dollár, P., Girshick, R.B.: Lvis: A dataset for large vocabulary instance segmentation. *CVPR* (2019)
- [29] Gupta, S., Davidson, J., Levine, S., Sukthankar, R., Malik, J.: Cognitive mapping and planning for visual navigation. *CVPR* pp. 2616–2625 (2017)
- [30] Hahn, M., Chaplot, D.S., Tulsiani, S.: No rl, no simulation: Learning to navigate without navigating. *NeurIPS* (2021)
- [31] Henry, P., Krainin, M., Herbst, E., Ren, X., Fox, D.: Rgb-d mapping: Using depth cameras for dense 3d modeling of indoor environments. *Experimental robotics* (2014)
- [32] Hu, P., Sclaroff, S., Saenko, K.: Uncertainty-aware learning for zero-shot semantic segmentation. *NeurIPS* (2020)
- [33] Jia, C., Yang, Y., Xia, Y., Chen, Y.T., Parekh, Z., Pham, H., Le, Q.V., Sung, Y., Li, Z., Duerig, T.: Scaling up visual and vision-language representation learning with noisy text supervision. *ICML* (2021)
- [34] Kadian, A., Truong, J., Gokaslan, A., Clegg, A., Wijmans, E., Lee, S., Savva, M., Chernova, S., Batra, D.: Are we making real progress in simulated environments? measuring the sim2real gap in embodied visual navigation. *arXiv* (2019)
- [35] Kato, N., Yamasaki, T., Aizawa, K.: Zero-shot semantic segmentation via variational mapping. *ICCV-W* (2019)
- [36] Khandelwal, A., Weihs, L., Mottaghi, R., Kembhavi, A.: Simple but effective: Clip embeddings for embodied ai. *arXiv* (2021)
- [37] Kuipers, B., Byun, Y.T.: A robot exploration and mapping strategy based on a semantic hierarchy of spatial representations. *Robotics and autonomous systems* (1991)
- [38] Li, Z., Yao, L., Zhang, X., Wang, X., Kanhere, S.S., Zhang, H.: Zero-shot object detection with textual descriptions. *AAAI* (2019)

- [39] Liang, Y., Chen, B., Song, S.: Sscnav: Confidence-aware semantic scene completion for visual semantic navigation. ICRA (2021)
- [40] Lin, T.Y., Maire, M., Belongie, S., Bourdev, L., Girshick, R., Hays, J., Perona, P., Ramanan, D., Zitnick, C.L., Dollár, P.: Microsoft coco: Common objects in context. ECCV (2014)
- [41] Mac, T.T., Copot, C., Tran, D.T., De Keyser, R.: Heuristic approaches in robot path planning: A survey. Robotics and Autonomous Systems (2016)
- [42] Mao, Q., Wang, C., Yu, S., Zheng, Y., Li, Y.: Zero-shot object detection with attributes-based category similarity. IEEE Transactions on Circuits and Systems II: Express Briefs (2020)
- [43] Mezghani, L., Sukhbaatar, S., Lavril, T., Maksymets, O., Batra, D., Bojanowski, P., Karteek, A.: Memory-augmented reinforcement learning for image-goal navigation. arXiv (2021)
- [44] Minderer, M., Djolonga, J., Romijnders, R., Hubis, F.A., Zhai, X., Houlsby, N., Tran, D., Lucic, M.: Revisiting the calibration of modern neural networks. arXiv (2021)
- [45] Montavon, G., Lapuschkin, S., Binder, A., Samek, W., Müller, K.R.: Explaining nonlinear classification decisions with deep taylor decomposition. Pattern recognition **65**, 211–222 (2017)
- [46] Moravec, H., Elfes, A.: High resolution maps from wide angle sonar. ICRA (1985)
- [47] Newcombe, R.A., Izadi, S., Hilliges, O., Molyneaux, D., Kim, D., Davison, A.J., Kohi, P., Shotton, J., Hodges, S., Fitzgibbon, A.: Kinectfusion: Real-time dense surface mapping and tracking. ISMAR (2011)
- [48] Olson, C.F., Matthies, L.H.: Maximum likelihood rover localization by matching range maps. ICRA (1998)
- [49] Parisi, S., Dean, V., Pathak, D., Gupta, A.K.: Interesting object, curious agent: Learning task-agnostic exploration. NeurIPS (2021)
- [50] Pathak, D., Agrawal, P., Efros, A.A., Darrell, T.: Curiosity-driven exploration by self-supervised prediction. ICML (2017)
- [51] Pham, H., Dai, Z., Ghiasi, G., Liu, H., Yu, A.W., Luong, M.T., Tan, M., Le, Q.V.: Combined scaling for zero-shot transfer learning. arXiv (2021)
- [52] Pratt, S., Weihs, L., Farhadi, A.: The introspective agent: Interdependence of strategy, physiology, and sensing for embodied agents. arXiv (2022)
- [53] Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning transferable visual models from natural language supervision. ICML (2021)
- [54] Rahman, S., Khan, S.H., Porikli, F.M.: Zero-shot object detection: Learning to simultaneously recognize and localize novel concepts. ACCV (2018)
- [55] Raileanu, R., Rocktäschel, T.: Ride: Rewarding impact-driven exploration for procedurally-generated environments. ICLR (2020)
- [56] Ramakrishnan, S.K., Chaplot, D.S., Al-Halah, Z., Malik, J., Grauman, K.: Poni: Potential functions for objectgoal navigation with interaction-free learning. arXiv (2022)
- [57] Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: Imagenet large scale visual recognition challenge. IJCV (2015)

- [58] Savva, M., Chang, A.X., Dosovitskiy, A., Funkhouser, T., Koltun, V.: Minos: Multimodal indoor simulator for navigation in complex environments. arXiv (2017)
- [59] Savva, M., Kadian, A., Maksymets, O., Zhao, Y., Wijmans, E., Jain, B., Straub, J., Liu, J., Koltun, V., Malik, J., Parikh, D., Batra, D.: Habitat: A Platform for Embodied AI Research. ICCV (2019)
- [60] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., Klimov, O.: Proximal policy optimization algorithms. arXiv (2017)
- [61] Selvaraju, R.R., Das, A., Vedantam, R., Cogswell, M., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. IJCV (2019)
- [62] Song, S., Zhang, L., Xiao, J.: Robot in a room: Toward perfect object recognition in closed environments. arXiv (2015)
- [63] Strehl, A.L., Littman, M.L.: An analysis of model-based interval estimation for markov decision processes. J. Comput. Syst. Sci. (2008)
- [64] Wani, S., Patel, S., Jain, U., Chang, A.X., Savva, M.: Multion: Benchmarking semantic map memory using multi-object navigation. NeurIPS (2020)
- [65] Weihs, L., Deitke, M., Kembhavi, A., Mottaghi, R.: Visual room rearrangement. CVPR (2021)
- [66] Weihs, L., Salvador, J., Kotar, K., Jain, U., Zeng, K.H., Mottaghi, R., Kembhavi, A.: Allenact: A framework for embodied ai research. arXiv (2020)
- [67] Wijmans, E., Kadian, A., Morcos, A.S., Lee, S., Essa, I., Parikh, D., Savva, M., Batra, D.: Dd-ppo: Learning near-perfect pointgoal navigators from 2.5 billion frames. ICLR (2019)
- [68] Wilcox, B.H.: Robotic vehicles for planetary exploration. Applied Intelligence (1992)
- [69] Wortsman, M., Ehsani, K., Rastegari, M., Farhadi, A., Mottaghi, R.: Learning to learn how to learn: Self-adaptive visual navigation using meta-learning. CVPR (2019)
- [70] Wortsman, M., Ilharco, G., Gadre, S.Y., Roelofs, R., Gontijo-Lopes, R., Morcos, A.S., Namkoong, H., Farhadi, A., Carmon, Y., Kornblith, S., Schmidt, L.: Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. arXiv (2022)
- [71] Wortsman, M., Ilharco, G., Li, M., Kim, J.W., Hajishirzi, H., Farhadi, A., Namkoong, H., Schmidt, L.: Robust fine-tuning of zero-shot models. arXiv (2021)
- [72] Xia, F., Zamir, A.R., He, Z., Sax, A., Malik, J., Savarese, S.: Gibson env: Real-world perception for embodied agents. CVPR (2018)
- [73] Yamauchi, B.: A frontier-based approach for autonomous exploration. CIRA (1997)
- [74] Yan, C., Misra, D., Bennnett, A., Walsman, A., Bisk, Y., Artzi, Y.: Chalet: Cornell house agent learning environment. arXiv (2018)
- [75] Yang, W., Wang, X., Farhadi, A., Gupta, A., Mottaghi, R.: Visual semantic navigation using scene priors. arXiv (2018)
- [76] Zhou, C., Loy, C.C., Dai, B.: Denseclip: Extract free dense labels from clip. arXiv (2021)
- [77] Zhu, P., Wang, H., Saligrama, V.: Zero shot detection. IEEE Transactions on Circuits and Systems for Video Technology (2020)
- [78] Zhu, Y., Mottaghi, R., Kolve, E., Lim, J.J., Gupta, A., Fei-Fei, L., Farhadi, A.: Target-driven visual navigation in indoor scenes using deep reinforcement learning. ICRA (2017)

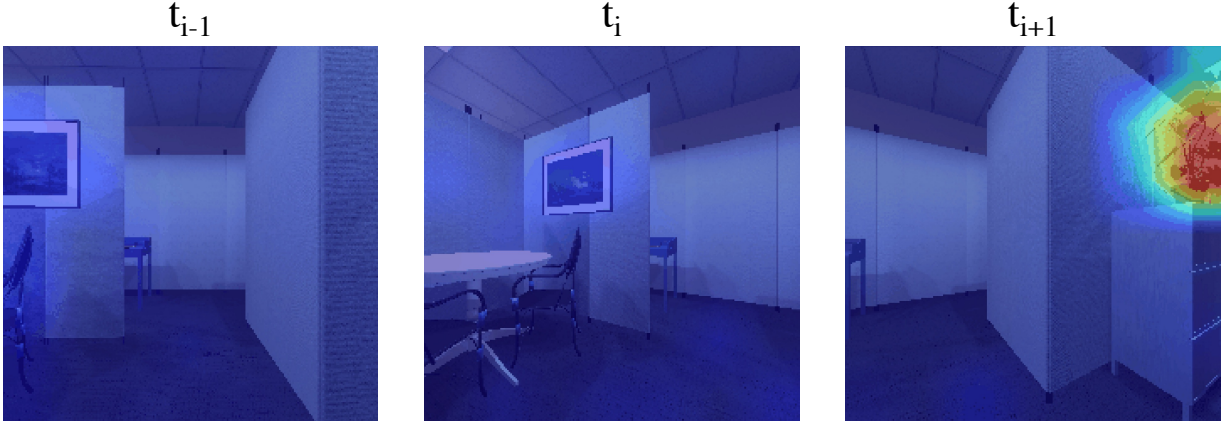


Figure 6: **Additional Grad-CAM localization visualization.** The target object is a plant. The localization is uniformly unconfident when the plant is not visible (frames t_{i-1} and t_i) and is confident when the plant is visible (frame t_{i+1}).

A Grad-CAM Details

Given a target class c , an input image x and a model f_θ , Grad-CAM [61] produces a localization map L_c , capturing the importance of each pixel for the image to be classified as the target class. For standard convolutional neural networks, neuron importance weights α_k^c are obtained from average pooling the gradients of the activations A^k :

$$\alpha_k^c = \text{AveragePool}_{i,j} \left(\frac{\partial y^c}{\partial A_{ij}^k} \right) \quad (1)$$

The saliency map is then given by a linear combination of the activations, weighted by the importance weights from Equation 1:

$$L_c = \text{ReLU} \left(\sum_k \alpha_k^c A^k \right), \quad (2)$$

where ReLU denotes the rectified linear unit function, $\text{ReLU}(x) = \max(0, x)$. The final saliency map is obtained by resizing L_c to the same size as the original image, using bilinear interpolation.

For ViTs [23], we follow the modifications of Chefer *et al.* [13]. Specifically, given the attention maps A^k for each transformer block k , and relevance scores R^k [45], the saliency map L_c^{ViT} is given by:

$$L_c^{\text{ViT}} = \prod_k \bar{A}^k, \quad (3)$$

$$\bar{A}^k = I + \mathbb{E}_h [\text{ReLU} (\nabla A^k \odot R^k)], \quad (4)$$

where \mathbb{E}_h computes the mean over the attention heads and \odot represents the element-wise product. In our case, we only look at attention maps from the final transformer block. Hence, $k = 1$.

In Fig. 6 we additionally visualize Grad-CAM saliency maps with “plant” as the target class. Notice that when the plant is not in view, saliency is qualitatively low. However, when the plant is in view, saliency spikes in the region of the plant. This visualization suggests that the localizer provides signal both for when the plant is not visible (the case where the agent should keep exploring) and when it is visible (the case where the agent should exploit this information to move to the plant).

Hyperparameter	Value
Discount factor (γ)	0.99
GAE parameter (λ)	0.95
Clipping parameter (ϵ)	0.1
Value loss coefficient	0.5
Entropy loss coefficient	0.01
LR	3e-4
Optimizer	Adam
Training steps	60M
Steps per rollout	500

Table 4: **DD-PPO hyperparameters.** Used for training exploration agents in both Habitat and RoboTHOR.

B Object Localization Details

We need a way to turn saliency maps from our localizers to determine regions of interest that correspond to the target. For example, if an agent is looking for a plant, a localizer should fire in the region of the plant. However, saliency maps—especially for Grad-CAM—are not directly interpretable as probabilities. Hence we empirically tune thresholds to extract localization information. For Grad-CAM we scale the ViT saliency maps L_c^{ViT} by 20 in Habitat and 10 in RoboTHOR. We then clip values between 0 and 1. If a pixel values exceeds 0.75, we consider that localizer to have “fired” for those pixels and register them in the voxel map as regions of interest (ROIs) for goal oriented navigation. The agent then plans towards such ROIs. For 9-Patch and 9-Lang, we have the advantage that we are softmax-ing over nine options in each case, thus creating a distribution over spatial regions. If confidence/probability exceeds 0.9 for a prediction, we consider the localizer to have “fired” for that patch.

C DD-PPO Hyperparameters

We use the AllenAct [66] framework to conduct exploration agent training. Our DD-PPO hyperparameters are in Tab. 4.

D Exploration Dataset Details

We adopt Habitat and RoboTHOR training scenes to fit our exploration agents. This set of rooms is disjoint from the downstream zero-shot object navigation testing rooms. Hence, rooms are unseen at test-time. In Habitat and RoboTHOR there are 61 and 60 train scenes respectively. For navigation training, starting positions are randomly chosen for each episode.

E Failure Checking Hyperparameters

Recall that we compare successive depth frames to determine if actions fail. We compute the depth difference $\Delta_{i,i+1} = |D_i - D_{i+1}|$. If mean $\mu(\Delta_{i,i+1}) < 0.1\text{m}$ and standard deviation $\sigma(\Delta_{i,i+1}) < 0.1\text{m}$, we consider the action to have failed.

F Object Localization Dataset Details

To create the object localization dataset, we sample agent poses around all possible goal objects across all categories and test rooms. Since there are less instances in RoboTHOR, we sample three poses for each

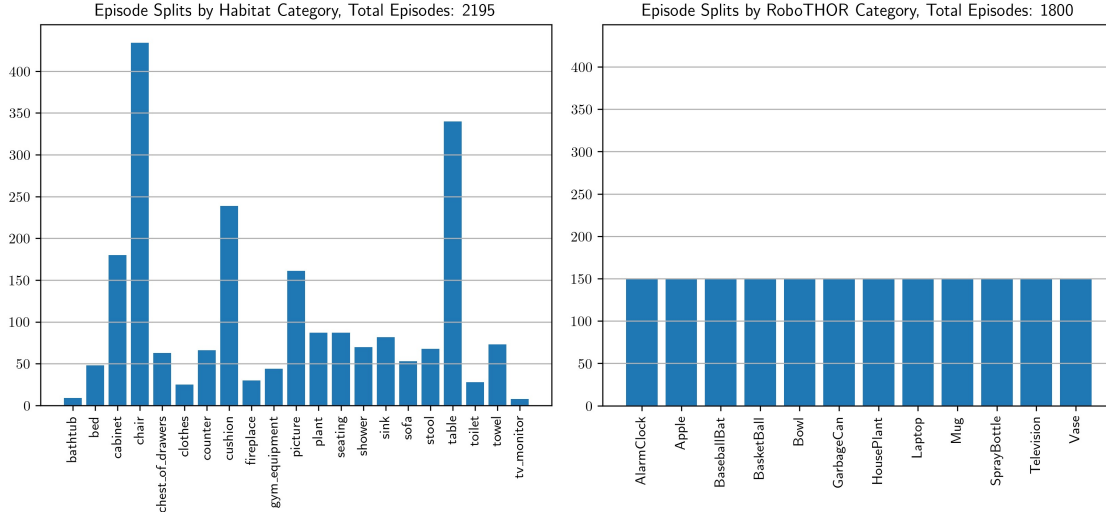


Figure 7: **Habitat and RoboTHOR episode splits.** Distribution of episodes in each CVPR 2021 object navigation challenge validation set that we adopt as our test set.

instance. In Habitat poses where at least one objects GT bounding box is larger than 0.4% of the image area and segmentation mask is larger than 0.2% of the image area are kept. Otherwise the pose is thrown out and another is sampled. In RoboTHOR we relax the box threshold to 0.09% as there are smaller objects in this dataset. Additionally we do not filter on masks.

G Object Navigation Dataset Details

We consider the following 21 categories in Habitat: chair, table, picture, cabinet, cushion, sofa, bed, chest_of_drawers, plant, sink, toilet, stool, towel, tv_monitor, shower, bathtub, counter, fireplace, gym_equipment, seating, clothes. We consider the following 12 categories in RoboTHOR: AlarmClock, Apple, BaseballBat, BasketBall, Bowl, GarbageCan, HousePlant, Laptop, Mug, SprayBottle, Television, Vase. Both of these lists include all objects for the Habitat and RoboTHOR CVPR 2021 object navigation challenges respectively. While the distribution of navigation episodes in RoboTHOR is equally split per category, this is not the case in Habitat object navigation validation split, which is our test split. We provide testing episode splits per category in Fig. 7.

H Additional Results

Results for category balancing in Habitat. Because of the uneven category-level split in the Habitat dataset seen in Fig. 7, we additionally provide the class balanced SPL and Success values in Tab. 5 (right). Class balancing is not done in the standard challenge Habitat metric. However, we feel it is important to note that the official metric is heavily biased to measure the performance on, for example, chair, cushion, and table categories.

Category-level performance. We also provide per category numbers for Habitat SPL (Tab. 6) and SUCCESS (Tab. 7), RoboTHOR SPL (Tab. 8) and SUCCESS (Tab. 9).

Exploration	CoW breeds		Fine-tuning steps	Habitat		Habitat (balanced)	
	Prompts	Localizer		SPL↑	SUCCESS↑	SPL↑	SUCCESS↑
FBE [73]	OpenAI	Grad-CAM [61]	0	0.063	0.111	0.036	0.064
FBE [73]	VG	Grad-CAM [61]	0	0.051	0.092	0.025	0.045
FBE [73]	OpenAI	9-Patch	0	0.047	0.088	0.025	0.051
FBE [73]	VG	9-Patch	0	0.042	0.084	0.025	0.052
Hab.-Learn	OpenAI	Grad-CAM [61]	60M	0.044	0.115	0.024	0.057
Hab.-Learn	VG	Grad-CAM [61]	60M	0.044	0.113	0.025	0.058
Robo.-Learn	OpenAI	Grad-CAM [61]	60M	0.032	0.086	0.020	0.046
Robo.-Learn	VG	Grad-CAM [61]	60M	0.029	0.078	0.018	0.042

Table 5: **Class balanced results in Habitat.** We provide class balanced results where we compute SPL and SUCCESS for each class individually and then average values over all classes. As can be seen by comparing the third and fourth columns, both metrics drop by about half. Note, for many classes there are very few samples, which yield noisier estimates. We again stress that we use the official validation split and (unbalanced) metric from the 2021 CVPR Habitat object navigation challenge in this paper.

category	FOG	FVG	FO9	FV9	HOG	HVG	ROG	RVG
bathtub	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
bed	0.013	0.000	0.011	0.000	0.019	0.000	0.000	0.000
cabinet	0.046	0.055	0.029	0.055	0.034	0.020	0.025	0.025
chair	0.042	0.032	0.055	0.041	0.045	0.046	0.038	0.031
chest_of_drawers	0.026	0.026	0.019	0.016	0.023	0.023	0.016	0.016
clothes	0.031	0.000	0.007	0.000	0.000	0.000	0.024	0.019
counter	0.040	0.025	0.031	0.001	0.009	0.022	0.010	0.009
cushion	0.058	0.038	0.059	0.018	0.069	0.055	0.039	0.037
fireplace	0.065	0.018	0.000	0.036	0.047	0.049	0.041	0.049
gym_equipment	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
picture	0.014	0.013	0.000	0.020	0.021	0.025	0.003	0.026
plant	0.031	0.014	0.022	0.012	0.024	0.019	0.049	0.023
seating	0.104	0.056	0.063	0.133	0.044	0.043	0.044	0.036
shower	0.021	0.033	0.029	0.022	0.026	0.041	0.029	0.021
sink	0.047	0.027	0.049	0.047	0.005	0.019	0.000	0.011
sofa	0.000	0.000	0.031	0.014	0.000	0.025	0.024	0.003
stool	0.005	0.007	0.008	0.008	0.000	0.018	0.005	0.000
table	0.207	0.185	0.117	0.106	0.107	0.117	0.069	0.063
toilet	0.000	0.000	0.000	0.000	0.017	0.000	0.000	0.000
towel	0.000	0.000	0.000	0.000	0.026	0.005	0.000	0.000
tv_monitor	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000

Table 6: **Category level SPL in Habitat.** F corresponds to FBE, H to Hab.-Learn, R to Robo.-Learn, O to OpenAI, V to VG, G to Grad-CAM, and 9 to 9-Patch.

category	FOG	FVG	FO9	FV9	HOG	HVG	ROG	RVG
bathtub	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
bed	0.021	0.000	0.021	0.000	0.021	0.000	0.000	0.000
cabinet	0.106	0.100	0.067	0.083	0.106	0.061	0.072	0.056
chair	0.071	0.053	0.101	0.062	0.090	0.092	0.071	0.062
chest_of_drawers	0.032	0.032	0.032	0.016	0.032	0.032	0.016	0.016
clothes	0.080	0.000	0.040	0.000	0.000	0.000	0.040	0.040
counter	0.091	0.045	0.061	0.015	0.015	0.045	0.030	0.015
cushion	0.092	0.063	0.092	0.025	0.113	0.096	0.075	0.059
fireplace	0.100	0.033	0.000	0.067	0.100	0.100	0.067	0.100
gym_equipment	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
picture	0.031	0.031	0.000	0.050	0.050	0.068	0.012	0.062
plant	0.046	0.023	0.046	0.023	0.046	0.034	0.069	0.034
seating	0.218	0.126	0.241	0.379	0.103	0.115	0.126	0.092
shower	0.029	0.043	0.029	0.043	0.043	0.071	0.043	0.043
sink	0.061	0.049	0.073	0.061	0.012	0.037	0.000	0.024
sofa	0.000	0.000	0.038	0.019	0.000	0.038	0.038	0.019
stool	0.015	0.015	0.015	0.015	0.000	0.044	0.015	0.000
table	0.356	0.332	0.212	0.232	0.382	0.376	0.282	0.256
toilet	0.000	0.000	0.000	0.000	0.036	0.000	0.000	0.000
towel	0.000	0.000	0.000	0.000	0.055	0.014	0.000	0.000
tv_monitor	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000

Table 7: **Category level Success in Habitat.** F corresponds to FBE, H to Hab.-Learn, R to Robo.-Learn, O to OpenAI, V to VG, G to Grad-CAM, and 9 to 9-Patch.

category	FOG	FVG	FO9	FV9	HOG	HVG	ROG	RVG
AlarmClock	0.041	0.046	0.031	0.021	0.022	0.021	0.061	0.040
Apple	0.078	0.095	0.097	0.093	0.101	0.095	0.095	0.085
BaseballBat	0.042	0.056	0.049	0.059	0.051	0.038	0.030	0.030
BasketBall	0.133	0.145	0.127	0.130	0.059	0.056	0.097	0.113
Bowl	0.034	0.036	0.048	0.049	0.020	0.004	0.036	0.031
GarbageCan	0.214	0.189	0.162	0.151	0.087	0.097	0.134	0.139
HousePlant	0.275	0.269	0.153	0.172	0.106	0.116	0.209	0.232
Laptop	0.117	0.089	0.086	0.095	0.050	0.056	0.080	0.065
Mug	0.042	0.029	0.104	0.118	0.085	0.053	0.055	0.062
SprayBottle	0.056	0.054	0.115	0.094	0.032	0.014	0.063	0.054
Television	0.134	0.098	0.115	0.075	0.082	0.028	0.109	0.075
Vase	0.039	0.035	0.028	0.039	0.026	0.040	0.036	0.039

Table 8: **Category level SPL in RoboTHOR.** F corresponds to FBE, H to Hab.-Learn, R to Robo.-Learn, O to OpenAI, V to VG, G to Grad-CAM, and 9 to 9-Patch.

category	FOG	FVG	FO9	FV9	HOG	HVG	ROG	RVG
AlarmClock	0.067	0.087	0.060	0.040	0.047	0.047	0.093	0.073
Apple	0.127	0.153	0.160	0.147	0.173	0.167	0.147	0.127
BaseballBat	0.100	0.128	0.093	0.100	0.133	0.127	0.080	0.080
BasketBall	0.160	0.167	0.167	0.167	0.113	0.113	0.107	0.140
Bowl	0.040	0.053	0.080	0.073	0.047	0.047	0.100	0.073
GarbageCan	0.333	0.273	0.287	0.233	0.193	0.220	0.233	0.207
HousePlant	0.413	0.396	0.240	0.260	0.253	0.267	0.320	0.333
Laptop	0.187	0.167	0.133	0.180	0.180	0.187	0.147	0.147
Mug	0.107	0.080	0.180	0.207	0.193	0.167	0.153	0.147
SprayBottle	0.080	0.093	0.187	0.160	0.093	0.047	0.107	0.107
Television	0.267	0.247	0.213	0.153	0.267	0.200	0.213	0.193
Vase	0.073	0.087	0.047	0.060	0.080	0.113	0.067	0.073

Table 9: **Category level Success in RoboTHOR.** F corresponds to FBE, H to Hab.-Learn, R to Robo.-Learn, O to OpenAI, V to VG, G to Grad-CAM, and 9 to 9-Patch.