# CoWs on PASTURE: Baselines and Benchmarks for Language-Driven Zero-Shot Object Navigation

Samir Yitzhak Gadre[◇]   Mitchell Wortsman[†]   Gabriel Ilharco[†]   Ludwig Schmidt[†]   Shuran Song[◇]

## Abstract

*For robots to be generally useful, they must be able to find arbitrary objects described by people (i.e., be language-driven) even without expensive navigation training on in-domain data (i.e., perform zero-shot inference). We explore these capabilities in a unified setting: language-driven zero-shot object navigation (L-ZSON). Inspired by the recent success of open-vocabulary models for image classification, we investigate a straightforward framework, CLIP on Wheels (CoW), to adapt open-vocabulary models to this task without fine-tuning. To better evaluate L-ZSON, we introduce the PASTURE benchmark, which considers finding uncommon objects, objects described by spatial and appearance attributes, and hidden objects described relative to visible objects. We conduct an in-depth empirical study by directly deploying 21 CoW baselines across HABITAT, ROBOTHOR, and PASTURE. In total, we evaluate over 90k navigation episodes and find that (1) CoW baselines often struggle to leverage language descriptions, but are proficient at finding uncommon objects. (2) A simple CoW, with CLIP-based object localization and classical exploration—and no additional training—matches the navigation efficiency of a state-of-the-art ZSON method trained for 500M steps on HABITAT MP3D data. This same CoW provides a 15.6 percentage point improvement in success over a state-of-the-art ROBOTHOR ZSON model.[1]*

## 1. Introduction

To be more widely applicable, robots should be language-driven: able to deduce goals based on arbitrary text input instead of being constrained to a fixed set of object categories. While existing image classification, semantic segmentation, and object navigation benchmarks like ImageNet-1k [61], ImageNet-21k [21], MS-COCO [43], LVIS [27], HABITAT [63], and ROBOTHOR [17] include a vast array of everyday items, they do not capture all objects that matter to people. For instance, a lost "toy airplane" may

◇Columbia University, †University of Washington. Correspondence to sy@cs.columbia.edu.

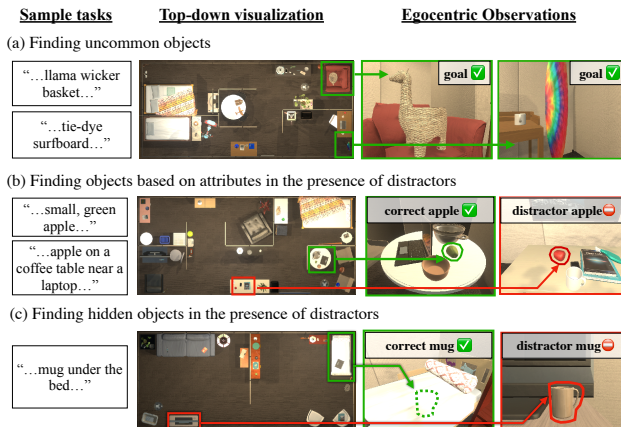[1]For code, data, and videos, see cow.cs.columbia.edu/



Figure 1. **The PASTURE benchmark for L-ZSON.** Text specifies navigation goal objects. Agents do not train on these tasks, making the evaluation protocol zero-shot. (a) *Uncommon object goals* like "llama wicker basket", not found in existing navigation benchmarks. (b) *Appearance, spatial descriptions*, which are necessary to find the correct object. (c) *Hidden object descriptions*, which localize objects that are not visible.

become relevant in a kindergarten classroom, but this object is not annotated in any of the aforementioned datasets.

In this paper, we study *Language-driven zero-shot object navigation (L-ZSON)*—a more challenging but also more applicable version of object navigation [4,17,63,74,83] and ZSON [37, 44] tasks. In L-ZSON, an agent must find an object based on a natural language description, which may contain different levels of granularity (e.g., "toy airplane", "toy airplane under the bed", or "wooden toy airplane"). L-ZSON encompasses ZSON, which specifies only the target category [37, 44]. Since L-ZSON is "zero-shot", we consider agents that do not have access to navigation training on the target objects or domains. This reflects realistic application scenarios, where the environment and object set may not be known a priori.

Performing L-ZSON in *any* environment with *unstructured* language input is challenging; however, recent advances in open-vocabulary models for image classification [34, 55, 57], object detection [3, 20, 26, 35, 41, 45, 47, 58, 82], and semantic segmentation [2, 5, 14, 32, 35, 36, 81] present a promising foundation. These models provide an

interface where one specifies—in text—the arbitrary objects they wish to classify, detect, or segment. For example, CLIP [57] open-vocabulary classifiers compute similarity scores between an input image and a set of user-specified captions (e.g., "a photo of a toy airplane.", ...), selecting the caption with the highest score to determine the image classification label. Given the flexibility of these models, we would like to understand their capability to execute embodied tasks even without additional training.

We present baselines and benchmarks for L-ZSON. More specifically:

- *A collection of baseline algorithms, CLIP on Wheels (CoW), which adapt open-vocabulary models to the task of L-ZSON.* CoW takes inspiration from the semantic mapping line of work [10, 40, 51], and decomposes the navigation task into exploration when the language target is not confidently localized, and target-driven planning otherwise. CoW retains the textual user interface of the original open-vocabulary model and does not require any navigation training. We evaluate 21 CoWs, ablating over many open-vocabulary models, exploration policies, backbones, prompting strategies, and post-processing strategies.

- *A new benchmark,* PASTURE*, to evaluate CoW and future methods on L-ZSON.* We design PASTURE, visualized in Fig. 1, to study capabilities that traditional object navigation agents, which are trained on a fixed set of categories, do not possess. We consider the ability to find (1) uncommon objects (e.g., "tie-dye surfboard"), (2) objects by their spatial and appearance attributes in the presence of distractor objects (e.g., "green apple" vs. "red apple"), and (3) objects that cannot be visually observed (e.g., "mug under the bed").

Together the CoW baselines and PASTURE benchmark allow us to conduct extensive studies on the capabilities of open-vocabulary models in the context of L-ZSON embodied tasks. Our experiments demonstrate CoW's potential on uncommon objects and limitations in taking full advantage of language descriptions—thereby providing empirical motivation for future studies. To contextualize CoW relative to prior zero-shot methods, we additionally evaluate on the HABITAT MP3D dataset. We find that our best CoW achieves navigation efficiency (SPL) that matches a state-of-the-art ZSON method [44] that trains on MP3D training data for 500M steps. On a ROBOTHOR object subset, considered in prior work, the same CoW beats the leading method [37] by 15.6 percentage points in task success.

## 2. Related Work

**Mapping and exploration.** Exploring effectively with a mobile robot is a long-standing problem in vision and robotics. Classical methods often decompose the task into map reconstruction [29, 31, 49, 50, 67], agent localization [16, 19, 52], and planning [40, 73]. Recent work investigates learned alternatives for exploration [6, 13, 53, 54, 59]. Here, agents are often trained end-to-end with self-supervised rewards (e.g., curiosity [54]) or supervised rewards (e.g., state visitation counts [24, 68, 70]). Learning-based methods typically need less hand-tuning, but require millions of training steps and reward engineering. We test both classical and learnable exploration strategies in the context of CoW to study their applicability to L-ZSON.

**Goal-conditioned navigation.** Apart from open-ended exploration, many navigation tasks are goal-conditioned, where the agent needs to navigate to a specified position (i.e., point goal [10, 11, 25, 30, 62, 72, 76, 78]), view of the environment (i.e., image goal [46, 60, 83]), or object category (i.e., object goal [1, 8, 9, 11, 18, 42, 69, 74, 79]). We consider an object goal navigation task.

**Instruction following in Navigation.** Prior work investigates language-based navigation, where language provides step-by-step instructions for the task [33, 38, 39]. This line of work demonstrates the benefits of additional language input for robot navigation, especially for long-horizon tasks (e.g., room-to-room navigation [39]). However, providing detailed step-by-step instructions (e.g., move 3 meters south [33]) could be challenging and time-consuming. In contrast, in our L-ZSON task, an algorithm gets natural language as the goal description instead of low-level instructions. Our task is more challenging as it requires the agent to infer its own exploration and searching strategies.

**Zero-shot object navigation (ZSON).** Recent work studies object navigation in zero-shot settings, where agents are evaluated on object categories that they are not explicitly trained on [37, 44]. Our task encompasses ZSON; however it also considers cases where more information—object attributes or hidden objects descriptions—is specified. Khandelwal *et al.* [37] train on a subset of ROBOTHOR categories and evaluate on a held-out set. In concurrent work, Majumdar *et al.* [44] train on an image goal navigation task and evaluate on object navigation downstream by leveraging CLIP multi-modal embeddings. Both algorithms necessitate navigation training for millions of steps and train separate models for each simulation domain. In contrast, CoW baselines do not necessitate any simulation training and can be deployed in multiple environments.

## 3. The L-ZSON Task

Language-driven zero-shot object navigation (L-ZSON) involves navigating to goal objects, specified in language, without explicit training to do so. Let $\mathcal{O}$ denote a set of natural language descriptions of target objects with potentially many attributes (e.g., "plant", "snake plant", "plant under the bed", etc.). This contrasts with definitions studied
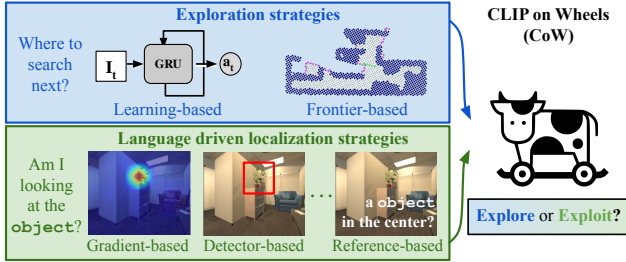
Figure 2. **CLIP on Wheels (CoW) overview.** A CoW uses a policy to explore and an object localizer (e.g., an open-vocabulary object detector) to determine if an object goal is in view.

in prior object navigation [4, 17] and ZSON [37, 44] tasks, which focus on high-level categories like "plant". Let $\mathcal{S}$ denote the set of navigation scenes. Let $p_0$ describe the initial pose of an agent. A navigation episode $\tau \in \mathcal{T}$ is written as a tuple $\tau = (s, o, p_0)$, $s \in \mathcal{S}, o \in \mathcal{O}$. Each $\tau$ is a *zero-shot task* as tuples of this form are not seen during training. Starting at $p_0$, an embodied agent's goal is to find $o$. The agent receives observations and sensor readings $I_t$ (e.g., RGB-D images). At each timestep $t$, the agent executes a navigation action $a \in \mathcal{A}$. A special action STOP $\in \mathcal{A}$ terminates the episode. If the agent is within $c$ units of $o$ and $o$ meets a visibility criteria, the episode is successful.

## 4. CLIP on Wheels (CoW) Baselines

To address L-ZSON, we investigate a simple baseline approach, CoW, which adapts open-vocabulary models like CLIP to make them suitable for the task. A CoW takes as input an egocentric RGB-D image and an object goal specified in language. As a CoW moves, it updates a top-down map of the world created using RGB-D observations and pose estimates (Sec. 4.1). Each CoW gets an exploration policy and a zero-shot object localization module as seen in Fig. 2. To observe diverse views of the scene, a CoW explores using a policy (Sec. 4.2). As the CoW roams, it keeps track of its confidence about the target object's location using an object localization module (Sec. 4.3) and its top-down map. When a CoW's confidence exceeds a threshold, it plans to the location of the goal and issues the STOP action. We now describe the ingredients used to make the CoWs evaluated in our experiments (Sec. 6).

### 4.1. Depth-based Mapping

As a CoW moves, it constructs a top-down map using input depth, pose estimates, and known agent height. The map is initialized using known camera intrinsics and the first depth image. Since a CoW knows the intended consequences of its actions (e.g., MOVEFORWARD should result in a 0.25m translation), each action is represented as a pose delta transform to approximate a transition. To deal with noise associated with actuation or depth, a CoW maintains a map at 0.125m resolution. To improve map accuracy,
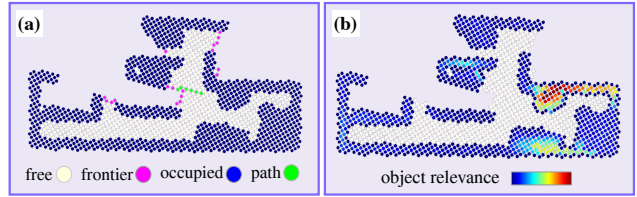


Figure 3. **Mapping.** Top-down map created from egocentric depth observations as a CoW roams a space. (a) Frontier Based Exploration [77] showing a planned path exploration path to the next frontier. (b) Back-projected object relevance scores provide object goal targets when a CoW has found an object.

a CoW checks for failed actions by comparing successive depth frames for movements (see Appx. A for details). Using known agent height (0.9m), map cells are projected to the ground plane to maintain a top-down representation of the world, which suffices for most navigation applications. Cells close to the floor are considered free space (white points in Fig. 3 (a)), while other cells are considered occupied (blue points in Fig. 3 (a)).

### 4.2. Exploration

Exploration generates diverse egocentric views so a CoW is more likely to view the language-specified target object. We consider two exploration methods, *frontier-based* and *learning-based*.

**Frontier based exploration (FBE) [77].** Using the top-down map discussed in Sec. 4.1, a CoW can navigate using a simple exploration heuristic: move to the frontier between free and unknown space to discover new regions. Once the navigator reaches a frontier (visualized as purple points in Fig. 3 (a)), it moves greedily to the next closest frontier. Since the map is updated at every timestep, a noisy pose estimate can contribute to inaccuracies. For example, narrow passages may collapse in the map due to pose drift. To circumvent such problems, we reinitialize the map when no paths exist to any frontiers in the map.

**Learnable exploration.** In addition to FBE, we investigate learnable alternatives, which may explore more intelligently. We consider an architecture and reward structure similar to prior work in embodied AI (e.g., [24, 37, 70]). Specifically, we adopt a frozen CLIP backbone with a trainable GRU [15] and linear heads for the actor and critic networks. We train agents independently in HABITAT [63] and ROBOTHOR [17] simulation environments for 60M steps each, using DD-PPO [65, 72] in the AllenAct [71] framework. We employ a simple count-based reward [68], which self-supervised exploration methods often attempt to approximate (e.g., [6]). All training scenes are disjoint from downstream navigation test scenes; however, learnable exploration involves millions of steps of simulation training, which the FBE heuristic does not necessitate. For details on reward, hyperparameters, and training, see Appx. B.

## 4.3. Object Localization

Successful navigation depends on object localization: the ability to tell *if* and *where* an object is in an image. Regions of high object relevance, extracted from 2D images, are projected to the depth-based map (Fig. 3 (b)) where they serve as natural navigation targets. To determine if and when a target is in an image, we consider the following object localization modules, used in our experiments (Sec. 6). For more details see Appx. C.

**CLIP [57] with $k$ referring expressions.** We match a CLIP visual embedding for the current observation against $k$ different CLIP text embeddings, which specify potential locations of the target object. For example, "a plant in the top left of the image." CLIP computes similarity between the image and text features to determine relevance scores over the language specified image regions.

**CLIP [57] with $k$ image patches.** We discretize the image into $k$ smaller patches, and run CLIP vision backbone inference on each of them to obtain patch features. We then convolve each patch feature with a CLIP text embedding for the target object. If the object is in a patch, the relevance score for that patch should be high.

**CLIP [57] with gradient relevance [12].** We adapt a network interpretability method [12, 66] to extract object relevancy from vision transformers (ViTs) [23]. Using a target CLIP text embedding and gradient information accumulated through the CLIP vision backbone, Chefer *et al*. [12] construct a relevance map over the image pixels, which qualitatively segments the target. These relevance methods typically assume the target object is visible and normalized between zero and one. We observe that removing this normalization produces high relevance when the target is in view and low relevance otherwise, hence providing signal for true positive and true negative detections.

**MDETR segmentation [35].** Kamath *et al*. [35] extend the DETR detector [7] to take arbitrary text and images as input and output box detections. They fine-tune their base model on PhraseCut [75], a dataset of paired masks and attribute descriptions, to support segmentation. We use this MDETR model to obtain object relevance for targets directly.

**OWL-ViT detection [47].** Minderer *et al*. [47] study a recipe to turn CLIP-like models into object detectors by fine-tuning on a set prediction task. Similar to MDETR, we use OWL-ViT to directly query images for targets.

**Post-processing.** The aforementioned models give continuous valued predictions. However, we are interested in binary masks giving if and where objects are in images. Hence, we threshold predictions for each model (see Appx. C for details). We further investigate two strategies for using the masks downstream: (1) using the entire mask or (2) using only the center pixel of masks. The second
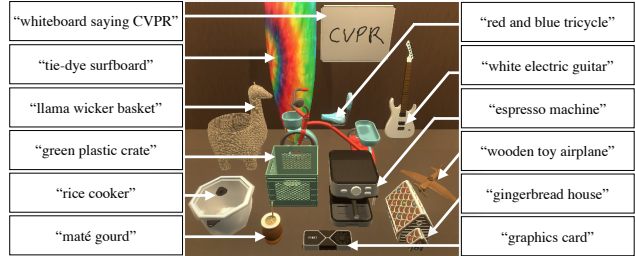


Figure 4. **Uncommon objects** in PASTURE.

strategy is reasonable because only part of an object needs to be detected for successful navigation.

**Target driven planning.** Recall, CoWs have depth sensors. We back-project object relevance from 2D images into the depth-based map (Sec. 4.1). We keep only the max relevance for each map cell (Fig. 3 (b)). CoWs can then plan to high relevance areas in the map. See Appx. D for additional method visualization.

## 5. The PASTURE Benchmark

To evaluate CoW baselines and future methods on L-ZSON, we introduce the PASTURE evaluation benchmark. PASTURE builds on ROBOTHOR validation scenes, which have parallel environments in the real-world. We target ROBOTHOR to facilitate future real-world benchmarking. PASTURE probes for seven capabilities explained in Sec. 5.1. We provide dataset statistics in Sec. 5.2.

### 5.1. PASTURE Tasks

PASTURE evaluates seven core L-ZSON capabilities.

**Uncommon objects.** Traditional benchmarks (e.g., ROBOTHOR and HABITAT MP3D) evaluate agents on common object categories like TVs; however, given the rich diversity of objects in homes, we would like to understand navigation performance on *uncommon objects*. Hence we add 12 new objects to each room. We use names shown in Fig. 4 as instance labels, which are minimal descriptions to identify each object. Some identifiers refer to text in images (e.g., "whiteboard saying CVPR") or to appearance attributes (e.g., "wooden toy airplane"). Other objects are less common in North America, like "maté", which is a popular Argentinian drink.

**Appearance descriptions.** To evaluate if baselines can take advantage of visual attributes, we introduce descriptions of the form: "{*size*}, {*color*}, {*material*} {*object*}". For example: "small, red apple", "orange basketball", "small, black, metallic alarm clock". Objects are considered small if their 3D bounding box diagonal is below a threshold. We determine color and materials by inspection.

**Spatial descriptions.** To test if agents can leverage spatial information in navigation, we introduce descriptions:

"{*object*} on top of {*x*}, near {*y*}, {*z*}, ...". For example, "house plant on a dresser near a spray bottle". To determine [on top of] relations, we use THOR metadata and to determine [nearness] we use a distance threshold between pairs of objects. We inspect all descriptions for correctness.

**Appearance descriptions with distractors.** To probe if appearance attributes better facilitate finding objects in the presence of distractors, we reuse the appearance captions from before, but evaluate on an modified environment with two visually distinct instances of each ROBOTHOR object category. For example, for the task of finding a "red apple", we have both a red apple and a green apple in the room. A navigator must leverage appearance information—and not just class information—to successfully complete the task.

**Spatial descriptions with distractors.** This capability is similar to the one above; however, we evaluate with spatial descriptions in the presence of distractor objects.

**Hidden object descriptions.** An ideal object navigator should find objects, even when they are hidden. To test this capability, we introduce descriptions: "{*object*} under/in {*x*}". For example, "basketball in the dresser drawers" or "vase under the sofa". We sample large objects (e.g., beds, sofas, dressers) in each scene to determine [under/in] relations. Additionally we remove visible instances of {*object*} from the room.

**Hidden object descriptions with distractors.** We use the hidden object descriptions from before, but reintroduce visible instances of {*object*} to serve as distractors. Consider finding a "mug under the bed". A distractor mug will also appear in the scene making the task more challenging.

## 5.2. Dataset Creation and Statistics

PASTURE contains three variations for each of the original 15 validation ROBOTHOR rooms: uncommon objects added, additional object instances added, and target objects removed. For each of the seven settings presented above, we evaluate over 12 object instances in 15 rooms with two initial agent starting locations. Hence PASTURE consists of 2,520 tasks, which is a similar order of magnitude to ROBOTHOR (1,800) and HABITAT MP3D (2,195) validation sets. For appearance attributes, 47% of the objects are considered "small". Each object gets an average of 1.2 color descriptors out of 22 possible choices, and 0.6 material descriptors out of 5 possible choices. Similarly, for spatial attributes, each object gets one object it is on top of or in (e.g., "vase in a shelving unit") and an average of 1.9 objects it is near. For a sample of appearance and spatial attributes see Fig. 1. For more dataset details and statistics see Appx. E.

## 6. Experiments

We first present our experimental setup, including the datasets, metrics, embodiment, and baselines considered

in our study (Sec. 6.1). Then we present results on PASTURE, thereby elucidating the strengths and weakness of CoW baselines for L-ZSON (Sec. 6.2). Finally, we compare to prior ZSON art in ROBOTHOR and HABITAT (MP3D) environments (Sec. 6.3).

### 6.1. Experimental setup

**Environments.** We consider PASTURE (Sec. 5), ROBOTHOR [17], and HABITAT (MP3D) [63] validation sets as our test sets. We utilize validation sets for testing because official test set ground-truth is not publicly available. Domains are setup with noise that is faithful to their original challenge settings. For ROBOTHOR—and by extension PASTURE—this means actuation noise but no depth noise. For HABITAT this means considerable depth noise and reconstruction artifacts, but no actuation noise.

**Navigation Metrics.** We adopt standard object navigation metrics to measure performance:

- SUCCESS (SR): the fraction of episodes where the agent executes STOP within 1.0m of the target object.
- Success weighted by inverse path length (SPL): Success weighted by the oracle shortest path length and normalized by the actual path length [4]. This metric points to the success efficiency of the agent.

In ROBOTHOR and PASTURE, the target must additionally be visible for the episode to be a success, which this is not the case in HABITAT—as specified in their 2021 challenge.

**Embodiment.** The agent is a LoCoBot [28]. All agents have discrete actions: {MOVEFORWARD, ROTATERIGHT, ROTATELEFT, STOP}. The move action advances the agent by 0.25m, while rotation actions pivot the camera by 30°.

**CoW Baselines.** For exploration we consider policies presented in Sec. 4.2: FBE heuristic exploration, learnable exploration optimized on HABITAT (MP3D) train scenes, and learned exploration optimized on ROBOTHOR train scenes. Learned exploration requires training in simulation—which is counter to our zero-shot goals; nonetheless, we ablate these explorers to contextualize their performance within the CoW framework. *FBE is the default CoW exploration strategy.*

For object localization, we consider:

- CLIP with $k = 9$ referring expressions (CLIP-Ref.)
- CLIP with $k = 9$ patches (CLIP-Patch)
- CLIP with gradient relevance (CLIP-Grad.)
- MDETR segmentation model (MDETR)
- OWL-ViT detector (OWL)

Descriptions of these models are in Sec. 4.3 and additional details are in Appx. C. All models are open-vocabulary. *No models are fine-tuned on navigation*, and hence we consider

| | CoW breeds | | | PASTURE | | | | | | | | | ROBOTHOR | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Uncom. | Appear. | Space | Appear. distract | Space distract | Hid. | Hid. distract | Avg. | | | |
| ID | Localizer | Arch. | Post | SR | SR | SR | SR | SR | SR | SR | SPL | SR | SPL | SR |
| ▲ | CLIP-Ref. | B/32 | | 2.8 | 1.4 | 1.4 | 0.8 | 1.4 | 4.7 | 5.0 | 1.2 | 2.5 | 1.6 | 2.2 |
| △ | CLIP-Ref. | B/32 | ✓ | 3.6 | 0.6 | 1.7 | 0.6 | 1.7 | 2.2 | 2.5 | 0.9 | 1.8 | 1.0 | 1.8 |
| ■ | CLIP-Ref. | B/16 | | 1.4 | 1.7 | 1.7 | 1.9 | 1.9 | 2.8 | 2.2 | 1.7 | 1.9 | 2.4 | 2.6 |
| □ | CLIP-Ref. | B/16 | ✓ | 1.4 | 2.8 | 2.8 | 3.1 | 3.3 | 1.7 | 1.9 | 1.7 | 2.4 | 2.1 | 2.7 |
| ▲ | CLIP-Patch | B/32 | | 10.6 | 9.7 | 6.7 | 6.4 | 6.4 | 16.7 | 16.7 | 7.5 | 10.4 | 9.0 | 14.3 |
| △ | CLIP-Patch | B/32 | ✓ | 18.1 | 13.3 | 13.3 | 8.6 | 10.8 | 17.5 | **17.8** | 9.0 | 14.2 | 10.6 | 20.3 |
| ■ | CLIP-Patch | B/16 | | 5.6 | 7.8 | 3.9 | 5.0 | 3.9 | 10.6 | 10.8 | 5.4 | 6.8 | 8.2 | 10.3 |
| □ | CLIP-Patch | B/16 | ✓ | 10.6 | 11.4 | 7.8 | 10.8 | 8.1 | 16.4 | 15.6 | 7.7 | 11.5 | 9.7 | 15.7 |
| ▲ | CLIP-Grad. | B/32 | | 13.6 | 10.6 | 9.2 | 7.5 | 7.2 | 13.9 | 12.8 | 8.3 | 10.7 | 9.6 | 13.8 |
| △ | CLIP-Grad. | B/32 | ✓ | 16.1 | 11.9 | 11.7 | 9.7 | 10.3 | 14.4 | 16.1 | 9.2 | 12.9 | 9.7 | 15.2 |
| ■ | CLIP-Grad. | B/16 | | 6.1 | 5.8 | 5.0 | 5.0 | 4.7 | 8.3 | 6.9 | 4.9 | 6.0 | 7.3 | 8.8 |
| □ | CLIP-Grad. | B/16 | ✓ | 8.1 | 10.8 | 8.6 | 8.6 | 6.7 | 11.1 | 11.4 | 6.7 | 9.3 | 8.6 | 11.6 |
| ◆ | MDETR | B3 | | 3.1 | 6.9 | 4.4 | 7.2 | 4.7 | 7.8 | 8.9 | 5.3 | 6.2 | 8.3 | 9.8 |
| ◇ | MDETR | B3 | ✓ | 3.1 | 7.2 | 5.0 | 6.9 | 4.7 | 8.1 | 8.9 | 5.4 | 6.3 | 8.4 | 9.9 |
| ▲ | OWL | B/32 | | 23.1 | 26.1 | 14.4 | 18.3 | 11.7 | 13.9 | 13.1 | 11.1 | 17.2 | 16.6 | 25.4 |
| △ | OWL | B/32 | ✓ | **32.8** | 26.4 | **19.4** | **19.4** | **16.1** | **19.2** | 14.4 | **12.6** | **21.1** | 16.9 | 26.7 |
| ■ | OWL | B/16 | | 25.8 | 23.6 | 15.3 | 17.2 | 12.5 | 13.1 | 13.9 | 11.4 | 17.3 | 16.2 | 24.8 |
| □ | OWL | B/16 | ✓ | 31.9 | **26.9** | 18.9 | **19.4** | 14.7 | 18.1 | 15.8 | **12.6** | 20.8 | **17.2** | **27.5** |
| | ProcTHOR fine-tune (supervised) [18] | | | n/a | n/a | n/a | n/a | n/a | n/a | n/a | n/a | n/a | 27.4 | 66.4 |

Table 1. **Benchmarking CoWs on PASTURE for L-ZSON.** On PASTURE we identify several key takeaways. (1) Average success on PASTURE is lower than on ROBOTHOR; however, CoWs are surprisingly good at finding uncommon objects (Uncom.), often finding them at higher rates than more common ROBOTHOR objects. (2) Comparing filled (●) vs. unfilled (○) row IDs, we notice post-processing mask predictions by using only the center pixel as a representative target for navigation helps in general (see Sec. 4.3 for more details on post-processing). (3) Comparing square (□) vs. triangle (△) IDs, we see that architectures (Arch.) using *more compute* (i.e., ViT-B/16) often perform comparably or worse than their competitors (i.e., ViT-B/32). This is especially true for CLIP [57] models (indicated in pink, orange, and purple). (4) Blue OWL-ViT [47] models perform best. (5) PASTURE tasks with distractor objects (distract) hurt performance and natural language specification is not sufficient to mitigate against the added difficulties in these tasks. (6) A supervised baseline shown in gray significantly outperforms CoWs on ROBOTHOR; however, it is unable to support PASTURE tasks out-of-the-box.

their inference *zero-shot* on our tasks.[2] We also consider various backbone architectures:

- A vision transformer [23], ViT-B/32 (▲ B/32)
- ViT-B/16 (■ B/16), which uses a smaller patch size of 16x16 and hence more compute.
- EfficentNet B3 (◆ B3), which is convolutional and similar in compute requirements to a ViT/B32.

For every model we additionally evaluate with post-processing (△, □, ◇), where only the center pixel of detections is registered in the top-down CoW map. Recall, this is a sensible strategy as only some part of the object needs to be found for an episode to be successful. For details on hyperparameters, learned agents, object localization threshold tuning, and CLIP prompt-tuning, see Appx. C, F.

**End-to-end learnable baselines.** We also compare against the following methods, which are trained in simulation for millions of steps:

- *EmbCLIP-ZSON* [37]: trains a model on eight ROBOTHOR categories, using CLIP language embeddings to specify the goal objects. At test time, the model is evaluated on four held-out object categories. The unseen target objects are specified with category names using CLIP language embeddings.
- *SemanticNav-ZSON* [44]: trains models separately—one for each dataset—for image goal navigation, where image goals are specified with CLIP visual embeddings. At test time image goals are switched with CLIP language embeddings for the object goals. We compare to the MP3D trained model.

Both EmbCLIP-ZSON and SemanticNav-ZSON leverage multi-modal CLIP visual and language embeddings in learnable frameworks that require simulation training.

### 6.2. CoWs on PASTURE

Tab. 1 shows the results of 18 CoWs evaluated on PASTURE. For category-level results see Appx. G. We now discuss several salient questions.

---

[2]The claim is not that these models have never seen any synthetic data in their large-scale training sets, only that they are not trained to navigate.
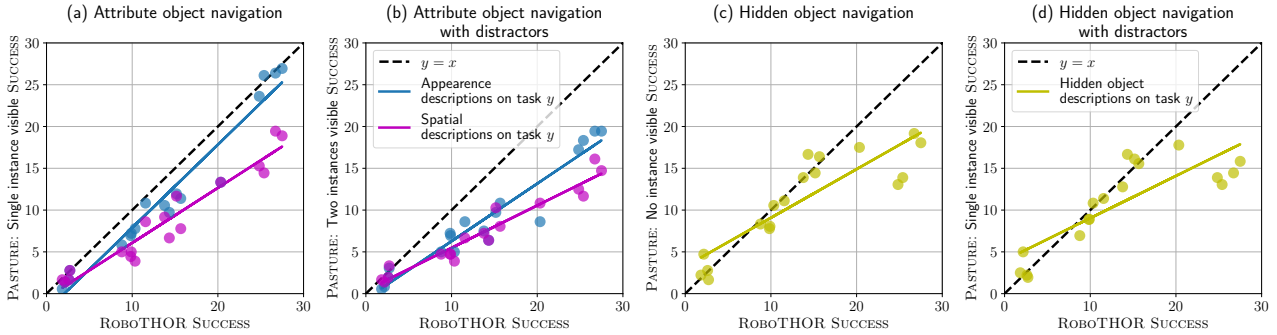
Figure 5. **PASTURE object navigation with descriptions.** In general object navigation with descriptions is more challenging than the ROBOTHOR object navigation task, as indicated by trend lines lying below the $y = x$ line. (a) Appearance descriptions are more helpful than spatial descriptions. (b) Performance further drops when distractor objects are introduced to the environment. However, CoWs are still able to make better use of appearance description than spatial descriptions. (c) Models in the lower success regime ($<15\%$ ROBOTHOR SUCCESS) perform comparably on finding hidden objects. However, this trend plateaus for higher success models. (d) Trends are similar when distractor objects are introduced for hidden object navigation.

**How well can CoWs find common objects vs. uncommon objects?** Comparing ROBOTHOR and uncommon (Uncom.) PASTURE success rate (SR) in Tab. 1—first and last columns—we notice that CoWs often find uncommon objects at higher rates than common ROBOTHOR objects (e.g., by ∼6 percentage points SUCCESS for the OWL ViT-B/32 CoW with post-post processing (△)). We hypothesize that though uncommon objects are less prevalent in daily life, they are still represented in open-vocabulary datasets and hence recognizable for the object localization modules. We further explore this hypothesis in Appx. E by visualizing CLIP retrieval results on LAION-5B [64] for the uncommon object categories. The relatively high performance on uncommon objects speaks to the flexibility of CoW baselines and their ability to inherit desirable properties from the open-vocabulary models.

**Can CoWs utilize appearance and spatial descriptions?** Looking at Fig. 5 (a) we see that neither appearance nor spatial descriptions improve CoW performance compared to their ROBOTHOR baseline performance (i.e., most points lie under the $y = x$ line). However, CoW is able to take better advantage of appearance descriptions than spatial descriptions. These results motivate future investigation on open-vocabulary object localization with a greater focus on textual object attributes.

**Can CoWs find visible objects in the presence of distractors?** In Fig. 5 (b) we see that CoWs experience further performance degradation when compared to Fig. 5 (a). We conclude that appearance and spatial attributes specified as natural language input are not sufficient to deal with the added complexity of distractors given current open-vocabulary models.

**Can CoWs find hidden objects?** Looking at Fig. 5 (c) we notice that models in the lower success regime (less that 15% SUCCESS on ROBOTHOR) are able to find hidden ob-

| | CoW breeds | | | | PASTURE (Avg.) | | ROBOTHOR | |
|---|---|---|---|---|---|---|---|---|
| ID | Loc. | Arch. | Post | Exp. Strategy | SPL | SR | SPL | SR |
| △ | OWL | B/32 | ✓ | ROBOTHOR learn. | 10.2 | 17.3 | 13.1 | 20.9 |
| △ | OWL | B/32 | ✓ | HABITAT learn. | 8.6 | 19.4 | 9.8 | 20.4 |
| △ | OWL | B/32 | ✓ | FBE | **12.6** | **21.1** | **16.9** | **26.7** |

Table 2. **Exploration ablation.** For a fixed object localizer (OWL-ViT B/32 with post processing), we ablate over different choices of exploration policy: the FBE heuristic, agents trained in ROBOTHOR, and HABITAT (MP3D). We find that FBE outperforms learnable alternatives on both PASTURE and ROBOTHOR. HABITAT learnable model perform worst, but are not trained on any PASTURE or ROBOTHOR data.

jects at about the same rate as ROBOTHOR objects (i.e., they lie on the $y = x$ line). OWL models in the higher success regime ($>15\%$) do not continue this trend; however, they do achieve higher absolute accuracy as seen in Tab. 1. Dealing with occlusion is a longstanding problem in computer vision, and these results provide a foundation upon which future hidden object navigation work can improve.

**Can CoWs find hidden objects in the presence of distractors?** Comparing Figs. 5 (c) and (d), we notice similar trends lines, with the best models performing worse with distractors. This suggests that distractors are less of a concern in the case of hidden objects than for visible object targets. In light of the fact that detection methods generally work better on larger objects, we hypothesize this effect is because distractor objects are smaller (e.g., apples, vases, basketballs) than objects used to conceal target categories (e.g., beds, sofas, etc.).

**What exploration method performs best?** We ablate the decision to use FBE for most experiments by fixing an object localizer (OWL, B/32, with post-processing (△)) and comparing against ROBOTHOR learnable exploration and HABITAT learnable exploration in Tab. 2.
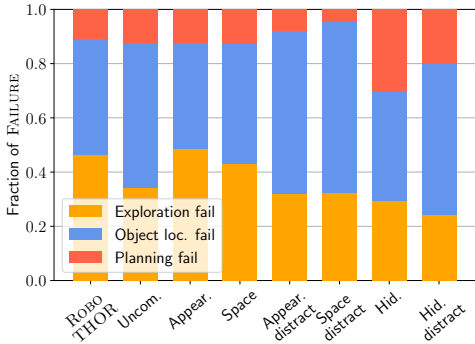
Figure 6. **Failure analysis for OWL, B/32, post-processing (△).** Exploration and object localization errors occur at similar ratios, with increased localization failures in the presense of distractors.

We notice that FBE performs best in all cases; however, learnable exploration still performs well suggesting that these models do learn useful strategies for the downstream tasks. Furthermore, the worse performance of the HABITAT learned model suggests that training on a certain domain can hurt in cases of inference in other domains (i.e., HABITAT ⟶ ROBOTHOR or HABITAT ⟶ PASTURE). Future work might re-consider learning-based algorithms for navigation that perform well out-of-distribution (i.e., "sim2real" and "sim2sim" exploration transfer as in [25]).

**How do CoWs fail?** We identify three high-level failure modes. (1) *Exploration fail*: the target is never seen. (2) *Object localization fail*: the target is seen but the localizer never fires. (3) *Planning fail*: the target is seen and the localizer fires, but planning fails due to inaccuracy in the map representation (Sec. 4.2). Looking at Fig. 6, we notice a large fraction of failures are due to exploration and object localization. This suggests CoWs may continue to improve as research in these fields progress. In Fig. 6 we additionally see that in cases where distractors are present a higher fraction of object localization failures occurs, supporting the claim that open-vocabulary models currently struggle to make full use of attribute prompts. See Appx. H for additional failure analysis.

### 6.3. Comparison to Prior Art

While our primary aim is to evaluate CoWs in general L-ZSON settings, we further evaluate CoWs in ZSON settings that prior work considers to establish CoW as a strong baseline for these tasks. Recall, ZSON can be seen as a special case of L-ZSON where only object goals are specified in language (no attributes).

In Tab. 3, we see there exists a CoW that outperforms the end-to-end baselines in all cases except SemanticNav-ZSON SUCCESS on HABITAT (MP3D). For instance, the CLIP-Grad., B/32, with post-processing (△) matches the SemanticNav-ZSON model on HABITAT (MP3D) SPL—

| ID | CoW breeds Loc. | Arch. | Post | HABITAT (MP3D) SPL | SR | ROBOTHOR (subset) SPL | SR | ROBOTHOR (full) SPL | SR | Nav. training steps |
|---|---|---|---|---|---|---|---|---|---|---|
| △ | CLIP-Grad. | B/32 | ✓ | **4.9** | 9.2 | 15.0 | 23.7 | 9.7 | 15.2 | **0** |
| △ | OWL | B/32 | ✓ | 3.7 | 7.4 | **20.8** | **32.5** | **16.9** | **26.7** | **0** |
| | EmbCLIP-ZSON [37] | | | – | – | – | 8.1 | – | 14.0* | 60M |
| | SemanticNav-ZSON [44] | | | 4.8 | **15.3** | – | – | – | – | 500M |

Table 3. **Comparison to prior art on existing ZSON benchmarks.** CoWs are able to match or out-compete existing methods that leverage millions of steps of navigation training in the evaluation simulator. *indicates a result from prior work that includes, non-zero-shot evaluation. Specifically, 1/4 of the evaluations are zero-shot on ROBOTHOR (subset) and the remaining 3/4 on categories seen during training.

4.9 for CoW v.s. 4.8 for the competitor, while also improving over EmbCLIP-ZSON ROBOTHOR by 15.6 percentage points. To contextualize the significance of this result, we reiterate that this CoW is trained for **0** navigation steps, while SemanticNav-ZSON and EmbCLIP-ZSON train in the target evaluation simulators for 500M and 60M steps respectively.

The superior performance of SemanticNav-ZSON in terms of MP3D SUCCESS indicates that there can be benefits to in-domain learning over CoW baselines. Future work may consider unifying the benefits of CoW-like models and fine-tuned models to get the best of both worlds.

## 7. Limitations and Conclusion

**Limitations.** While our evaluation of CoWs on HABITAT, ROBOTHOR, and PASTURE is a step towards assessing their performance in different domains, ultimately, real-world performance matters most. Hence, the biggest limitation of our work is the lack of large-scale, real-world benchmarking—which is also missing in much of the related literature. Additionally, CoW inherents the meta-limitations of the object localization and exploration methods considered. For example, object localizers require tuning a confidence threshold to balance precision and recall. Finally, we do not consider different agent embodiment or continuous action spaces. This is a pertinent investigation given recent findings of Pratt *et al.* [56] that agent morphology can be a big determinant of downstream performance.

**Conclusion.** This paper introduces the PASTURE benchmark for language-driven zero-shot object navigation and several CLIP on Wheels baselines, translating the successes of existing zero-shot models to an embodied task. We view CoW as an instance of using open-vocabulary models, with text-based interfaces, to tackle robotics tasks in more flexible settings. We hope that the baselines and the proposed benchmark will spur the field to explore broader and more powerful forms of zero-shot embodied AI.

# References

[1] Ziad Al-Halah, Santhosh K. Ramakrishnan, and Kristen Grauman. Zero experience required: Plug & play modular transfer learning for semantic visual navigation. *arXiv*, 2022. 2

[2] Donghyeon Baek, Youngmin Oh, and Bumsub Ham. Exploiting a joint embedding space for generalized zero-shot semantic segmentation. *ICCV*, 2021. 1

[3] Ankan Bansal, Karan Sikka, Gaurav Sharma, Rama Chellappa, and Ajay Divakaran. Zero-shot object detection. *ECCV*, 2018. 1

[4] Dhruv Batra, Aaron Gokaslan, Aniruddha Kembhavi, O. Maksymets, R. Mottaghi, M. Savva, A. Toshev, and Erik Wijmans. Objectnav revisited: On evaluation of embodied agents navigating to objects. *arXiv*, abs/2006.13171, 2020. 1, 3, 5

[5] Maxime Bucher, Tuan-Hung VU, Matthieu Cord, and Patrick Pérez. Zero-shot semantic segmentation. *NeurIPS*, 2019. 1

[6] Yuri Burda, Harrison Edwards, Amos J. Storkey, and Oleg Klimov. Exploration by random network distillation. *ICLR*, 2018. 2, 3

[7] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. *ECCV*, 2020. 4, 14

[8] Matthew Chang, Arjun Gupta, and Saurabh Gupta. Semantic visual navigation by watching youtube videos. *NeurIPS*, 2020. 2

[9] Devendra Singh Chaplot, Dhiraj Gandhi, Abhinav Gupta, and Ruslan Salakhutdinov. Object goal navigation using goal-oriented semantic exploration. *NeurIPS*, 2020. 2

[10] Devendra Singh Chaplot, Dhiraj Gandhi, Saurabh Gupta, Abhinav Kumar Gupta, and Ruslan Salakhutdinov. Learning to explore using active neural slam. *ICLR*, 2020. 2

[11] Prithvijit Chattopadhyay, Judy Hoffman, Roozbeh Mottaghi, and Aniruddha Kembhavi. Robustnav: Towards benchmarking robustness in embodied navigation. *ICCV*, 2021. 2

[12] Hila Chefer, Shir Gur, and Lior Wolf. Transformer interpretability beyond attention visualization. *CVPR*, 2021. 4, 13, 14

[13] Tao Chen, Saurabh Gupta, and Abhinav Gupta. Learning exploration policies for navigation. *ICLR*, 2019. 2

[14] Jiaxin Cheng, Soumyaroop Nandi, Prem Natarajan, and Wael Abd-Almageed. Sign: Spatial-information incorporated generative network for generalized zero-shot semantic segmentation. *ICCV*, 2021. 1

[15] Kyunghyun Cho, Bart van Merrienboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of neural machine translation: Encoder–decoder approaches. In *SSST@EMNLP*, 2014. 3

[16] Andrew J Davison and David W Murray. Mobile robot localisation using active vision. *ECCV*, 1998. 2

[17] Matt Deitke, Winson Han, Alvaro Herrasti, Aniruddha Kembhavi, Eric Kolve, Roozbeh Mottaghi, Jordi Salvador, Dustin Schwenk, Eli VanderBilt, Matthew Wallingford, Luca Weihs, Mark Yatskar, and Ali Farhadi. Robothor: An open simulation-to-real embodied ai platform. *CVPR*, 2020. 1, 3, 5

[18] Matt Deitke, Eli VanderBilt, Alvaro Herrasti, Luca Weihs, Jordi Salvador, Kiana Ehsani, Winson Han, Eric Kolve, Ali Farhadi, Aniruddha Kembhavi, and Roozbeh Mottaghi. Procthor: Large-scale embodied ai using procedural generation. *NeurIPS*, 2022. 2, 6

[19] Frank Dellaert, Dieter Fox, Wolfram Burgard, and Sebastian Thrun. Monte carlo localization for mobile robots. *ICRA*, 1999. 2

[20] Berkan Demirel, Ramazan Gokberk Cinbis, and Nazli Ikizler-Cinbis. Zero-shot object detection by hybrid region embedding. *BMVC*, 2018. 1

[21] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, K. Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. *CVPR*, 2009. 1

[22] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv*, 2020. 13

[23] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2020. 4, 6

[24] Samir Yitzhak Gadre, Kiana Ehsani, Shuran Song, and Roozbeh Mottaghi. Continuous scene representations for embodied ai. *CVPR*, 2022. 2, 3

[25] Daniel Gordon, Abhishek Kadian, Devi Parikh, Judy Hoffman, and Dhruv Batra. Splitnet: Sim2sim and task2task transfer for embodied visual navigation. *CVPR*, 2019. 2, 8

[26] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Zero-shot detection via vision and language knowledge distillation. *arXiv*, 2021. 1

[27] Agrim Gupta, Piotr Dollár, and Ross B. Girshick. Lvis: A dataset for large vocabulary instance segmentation. *CVPR*, 2019. 1

[28] Abhinav Kumar Gupta, Adithyavairavan Murali, Dhiraj Gandhi, and Lerrel Pinto. Robot learning in homes: Improving generalization and reducing dataset bias. *NeurIPS*, 2018. 5

[29] Saurabh Gupta, James Davidson, Sergey Levine, Rahul Sukthankar, and Jitendra Malik. Cognitive mapping and planning for visual navigation. *CVPR*, pages 2616–2625, 2017. 2

[30] Meera Hahn, Devendra Singh Chaplot, and Shubham Tulsiani. No rl, no simulation: Learning to navigate without navigating. *NeurIPS*, 2021. 2

[31] Peter Henry, Michael Krainin, Evan Herbst, Xiaofeng Ren, and Dieter Fox. Rgb-d mapping: Using depth cameras for dense 3d modeling of indoor environments. *Experimental robotics*, 2014. 2

[32] Ping Hu, Stan Sclaroff, and Kate Saenko. Uncertainty-aware learning for zero-shot semantic segmentation. *NeurIPS*, 2020. 1

[33] Chenguang Huang, Oier Mees, Andy Zeng, and Wolfram Burgard. Visual language maps for robot navigation. *arXiv preprint arXiv:2210.05714*, 2022. 2

[34] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V Le, Yunhsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. *ICML*, 2021. 1, 13

[35] Aishwarya Kamath, Mannat Singh, Yann LeCun, Ishan Misra, Gabriel Synnaeve, and Nicolas Carion. Mdetr - modulated detection for end-to-end multi-modal understanding. *ICCV*, 2021. 1, 4, 13, 14

[36] Naoki Kato, Toshihiko Yamasaki, and Kiyoharu Aizawa. Zero-shot semantic segmentation via variational mapping. *ICCV-W*, 2019. 1

[37] Apoorv Khandelwal, Luca Weihs, Roozbeh Mottaghi, and Aniruddha Kembhavi. Simple but effective: Clip embeddings for embodied ai. *arXiv*, 2021. 1, 2, 3, 6, 8

[38] Jacob Krantz, Erik Wijmans, Arjun Majumdar, Dhruv Batra, and Stefan Lee. Beyond the nav-graph: Vision-and-language navigation in continuous environments. In *European Conference on Computer Vision*, pages 104–120. Springer, 2020. 2

[39] Alexander Ku, Peter Anderson, Roma Patel, Eugene Ie, and Jason Baldridge. Room-across-room: Multilingual vision-and-language navigation with dense spatiotemporal grounding. *arXiv preprint arXiv:2010.07954*, 2020. 2

[40] Benjamin Kuipers and Yung-Tai Byun. A robot exploration and mapping strategy based on a semantic hierarchy of spatial representations. *Robotics and autonomous systems*, 1991. 2

[41] Zhihui Li, Lina Yao, Xiaoqin Zhang, Xianzhi Wang, Salil S. Kanhere, and Huaxiang Zhang. Zero-shot object detection with textual descriptions. *AAAI*, 2019. 1

[42] Yiqing Liang, Boyuan Chen, and Shuran Song. Sscnav: Confidence-aware semantic scene completion for visual semantic navigation. *ICRA*, 2021. 2

[43] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context. *ECCV*, 2014. 1

[44] Arjun Majumdar, Gunjan Aggarwal, Bhavika Devnani, Judy Hoffman, and Dhruv Batra. Zson: Zero-shot object-goal navigation using multimodal goal embeddings. *arXiv preprint arXiv:2206.12403*, 2022. 1, 2, 3, 6, 8

[45] Qiaomei Mao, Chong Wang, Sheng Yu, Ye Zheng, and Yuqi Li. Zero-shot object detection with attributes-based category similarity. *IEEE Transactions on Circuits and Systems II: Express Briefs*, 2020. 1

[46] Lina Mezghani, Sainbayar Sukhbaatar, Thibaut Lavril, Oleksandr Maksymets, Dhruv Batra, Piotr Bojanowski, and Alahari Karteek. Memory-augmented reinforcement learning for image-goal navigation. *arXiv*, 2021. 2

[47] Matthias Minderer, Alexey A. Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, Xiao Wang, Xiaohua Zhai, Thomas Kipf, and Neil Houlsby. Simple open-vocabulary object detection with vision transformers. *ECCV*, 2022. 1, 4, 6, 13, 14

[48] Grégoire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek, and Klaus-Robert Müller. Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern recognition*, 65:211–222, 2017. 13

[49] Hans Moravec and Alberto Elfes. High resolution maps from wide angle sonar. *ICRA*, 1985. 2

[50] Richard A Newcombe, Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J Davison, Pushmeet Kohi, Jamie Shotton, Steve Hodges, and Andrew Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. *ISMAR*, 2011. 2

[51] Lachlan Nicholson, Michael Milford, and Niko Sünderhauf. Quadricslam: Dual quadrics from object detections as landmarks in object-oriented slam. *RA-L*, 2019. 2

[52] Clark F Olson and Larry H Matthies. Maximum likelihood rover localization by matching range maps. *ICRA*, 1998. 2

[53] Simone Parisi, Victoria Dean, Deepak Pathak, and Abhinav Kumar Gupta. Interesting object, curious agent: Learning task-agnostic exploration. *NeurIPS*, 2021. 2

[54] Deepak Pathak, Pulkit Agrawal, Alexei A. Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. *ICML*, 2017. 2

[55] Hieu Pham, Zihang Dai, Golnaz Ghiasi, Hanxiao Liu, Adams Wei Yu, Minh-Thang Luong, Mingxing Tan, and Quoc V. Le. Combined scaling for zero-shot transfer learning. *arXiv*, 2021. 1, 13

[56] Sarah Pratt, Luca Weihs, and Ali Farhadi. The introspective agent: Interdependence of strategy, physiology, and sensing for embodied agents. *arXiv*, 2022. 8

[57] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. *ICML*, 2021. 1, 2, 4, 6, 13

[58] Shafin Rahman, Salman Hameed Khan, and Fatih Murat Porikli. Zero-shot object detection: Learning to simultaneously recognize and localize novel concepts. *ACCV*, 2018. 1

[59] Roberta Raileanu and Tim Rocktäschel. Ride: Rewarding impact-driven exploration for procedurally-generated environments. *ICLR*, 2020. 2

[60] Santhosh K. Ramakrishnan, Devendra Singh Chaplot, Ziad Al-Halah, Jitendra Malik, and Kristen Grauman. Poni: Potential functions for objectgoal navigation with interaction-free learning. *CVPR*, 2022. 2

[61] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large scale visual recognition challenge. *IJCV*, 2015. 1

[62] Manolis Savva, Angel X Chang, Alexey Dosovitskiy, Thomas Funkhouser, and Vladlen Koltun. Minos: Multimodal indoor simulator for navigation in complex environments. *arXiv*, 2017. 2

[63] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv Batra. Habitat: A Platform for Embodied AI Research. *ICCV*, 2019. 1, 3, 5

[64] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *NeurIPS*, 2022. 7, 16, 18

[65] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv*, 2017. 3

[66] Ramprasaath R. Selvaraju, Abhishek Das, Ramakrishna Vedantam, Michael Cogswell, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. *IJCV*, 2019. 4, 13, 14

[67] Shuran Song, Linguang Zhang, and Jianxiong Xiao. Robot in a room: Toward perfect object recognition in closed environments. *arXiv*, 2015. 2

[68] Alexander L. Strehl and Michael L. Littman. An analysis of model-based interval estimation for markov decision processes. *J. Comput. Syst. Sci.*, 2008. 2, 3

[69] Saim Wani, Shivansh Patel, Unnat Jain, Angel X. Chang, and Manolis Savva. Multion: Benchmarking semantic map memory using multi-object navigation. *NeurIPS*, 2020. 2

[70] Luca Weihs, Matt Deitke, Aniruddha Kembhavi, and Roozbeh Mottaghi. Visual room rearrangement. *CVPR*, 2021. 2, 3

[71] Luca Weihs, Jordi Salvador, Klemen Kotar, Unnat Jain, Kuo-Hao Zeng, Roozbeh Mottaghi, and Aniruddha Kembhavi. Allenact: A framework for embodied ai research. *arXiv*, 2020. 3, 12

[72] Erik Wijmans, Abhishek Kadian, Ari S. Morcos, Stefan Lee, Irfan Essa, D. Parikh, Manolis Savva, and Dhruv Batra. Ddppo: Learning near-perfect pointgoal navigators from 2.5 billion frames. *ICLR*, 2019. 2, 3

[73] Brian H Wilcox. Robotic vehicles for planetary exploration. *Applied Intelligence*, 1992. 2

[74] Mitchell Wortsman, Kiana Ehsani, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Learning to learn how to learn: Self-adaptive visual navigation using meta-learning. *CVPR*, 2019. 1, 2

[75] Chenyun Wu, Zhe Lin, Scott D. Cohen, Trung Bui, and Subhransu Maji. Phrasecut: Language-based image segmentation in the wild. *CVPR*, 2020. 4, 14

[76] Fei Xia, Amir R Zamir, Zhiyang He, Alexander Sax, Jitendra Malik, and Silvio Savarese. Gibson env: Real-world perception for embodied agents. *CVPR*, 2018. 2

[77] Brian Yamauchi. A frontier-based approach for autonomous exploration. *CIRA*, 1997. 3

[78] Claudia Yan, Dipendra Misra, Andrew Bennnett, Aaron Walsman, Yonatan Bisk, and Yoav Artzi. Chalet: Cornell house agent learning environment. *arXiv*, 2018. 2

[79] Wei Yang, Xiaolong Wang, Ali Farhadi, Abhinav Gupta, and Roozbeh Mottaghi. Visual semantic navigation using scene priors. *arXiv*, 2018. 2

[80] Xiaohua Zhai, Xiao Wang, Basil Mustafa, Andreas Steiner, Daniel Keysers, Alexander Kolesnikov, and Lucas Beyer. Lit: Zero-shot transfer with locked-image text tuning. *CVPR*, 2022. 14

[81] Chong Zhou, Chen Change Loy, and Bo Dai. Denseclip: Extract free dense labels from clip. *arXiv*, 2021. 1

[82] Pengkai Zhu, Hanxiao Wang, and Venkatesh Saligrama. Zero shot detection. *IEEE Transactions on Circuits and Systems for Video Technology*, 2020. 1

[83] Yuke Zhu, Roozbeh Mottaghi, Eric Kolve, Joseph J Lim, Abhinav Gupta, Li Fei-Fei, and Ali Farhadi. Target-driven visual navigation in indoor scenes using deep reinforcement learning. *ICRA*, 2017. 1, 2

# Contents

## A. Depth-based Mapping Details

**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., MOVEA-HEAD), which could be executed repeatedly.

To improve the action success of our CoWs, 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 $\mu$, $\sigma$ below which we can be confident actions failed. In our studies, we set these to $\mu = 0.1$m, $\sigma = 0.1$m.

Note, that in modern robot navigation systems, on-board pose estimation and bumper sensors for obstacle avoidance are nearly ubiquitous. Hence, it is possible to implement action failure checking beyond visual inputs, even though in this paper *we only consider vision-based failure checking*.

**Formalization of registering new depth observations.** Recall, a CoW constructs a top-down map based on ego-centric depth observation and approximated poses deltas. We provide a formalization of this process.

To create this map, we first estimate the pose of the CoW's coordinate frame. Let $C_i$ be the current local coordinate frame of a CoW at timestep $i$. Let $W$ denote a world frame. We would like to align observations from each local frame to $W$ to keep a single, consistent map. When a CoW is initialized, we set $^W\hat{T}_{C_0} = I$, where $I$ is the identity SE(3) transform and $^W\hat{T}_{C_0}$ is a transform, which aligns frame $C_0$ to $W$. 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 delta transform, which models an action transition. By concatenating these transforms over time, we get pose estimates. To estimate 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 action failures, these estimates are sensitive to pose drift. We use the current estimated pose $^W T_{C_t}$, camera intrinsic matrix $K$, current depth observation $D_t$, standard back-projection of $D_t$, and pose concatenation to register new observations in 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 down gravity direction. Points already near the floor are considered free space, while other points are considered occupied as shown in. Additionally, we discretize the map, which additionally helps deal with sensor noise.

## B. Exploration Details

**Learnable exploration.** We use the AllenAct [71] framework to conduct exploration agent training. Our DD-PPO hyperparameters are in Tab. 4. We employ a simple state visitation count based reward. 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 train two agents, one on ROBOTHOR and the other on HABITAT MP3D train sets, which are disjoint from the downstream validation sets we use for testing.

| Hyperparameter | Value |
|---|---|
| Discount factor ($\gamma$) | 0.99 |
| GAE parameter ($\lambda$) | 0.95 |
| Clipping parameter ($\epsilon$) | 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.

## C. Localization Details

Here we provide a review of the CLIP model [57], the concept of prompt ensembling as it relates to CLIP, the gradient saliency method introduced by Chefer *et al.* [12], MDETR [35], and OWL-ViT [47].

**CLIP overview.** Recent open-vocabulary models for zero-shot image classification include CLIP [57], ALIGN [34], and BASIC [55]. These methods pre-train contrastively on image-caption pairs from the internet. In the case of the original CLIP model, on 400M pairs. The key insight is that the image representation—extracted by a vision tower—and the caption representation—extracted by a text tower—should be similar. Hence, the contrastive objects encourages image and text representations for a positive pairs to be similar and for these representations to be dissimilar from other images and captions.

More formally, CLIP jointly trains two encoders: a visual encoder $f_\theta$ and a language encoder $f_\lambda$. Given a dataset of image-text pairs $\mathcal{D} = \{..., (I_k, t_k), (I_{k+1}, t_{k+1}), ...\}$, which can be generated at massive scale by leveraging internet data. CLIP employs standard mini-batch style training. Given a batch of size $n$, with $n$ image-text pairs, the current functions $f_\theta$ and $f_\lambda$ are used to featurize the data, yielding *L-2 normalized* embeddings for the batch: $\{(z_0^I, z_0^t), (z_1^I, z_1^t), ...(z_n^I, z_n^t), \}$. It is then possible to define a symmetric contrastive loss of the following form, for each sample in the batch:

$$\mathcal{L}(z_i^I, z_i^t) = \frac{1}{2}\left( \frac{z_i^I \cdot z_i^t}{\sum_{j \neq i} z_i^I \cdot z_j^t} + \frac{z_i^I \cdot z_i^t}{\sum_{j \neq i} z_j^I \cdot z_i^t} \right). \quad (1)$$

By minimizing this loss, the representation for images and their corresponding captions are being pushed closer together, while these representations are being pushed farther from other images and text features in the batch, which are assumed to contain dissimilar concepts.

After training, these models can be thought to match images and captions. Given a set of captions identifying concepts (e.g., "a photo of an apple.", "a photo of an orange", etc.) it is possible to construct a classification head by extracting text features, via $f_\lambda$, and L-2 normalizing each of them. Now given an image, say of an apple, we can extract a visual feature, via $f_\theta$, and again L-2 normalize. By dotting the visual feature with the text features, we can find the highest similarity score amongst the text captions and assign the image the corresponding label (e.g., "apple" from the caption "a photo of an apple."). Because this downstream classifier has not been trained to specifically classify images (say apples), it is considered a zero-shot classifier.

**CLIP prompt ensembling.** Radford *et al.* [57] find a simple strategy to boost performance of a CLIP zero-shot classifier. Instead of using one prompt for a class (e.g., "a photo of a apple."), they instead compute many text features for a class (e.g., "a photo of a apple.", "a blurry photo of a apple.", etc.). By simply averaging all these text features, they find performance improves. For CLIP-Ref., CLIP-Patch, and CLIP-Grad. strategies we similarly use prompt ensembling using the set of 80 prompt templates used for ImageNet-1k evaluation.[3] We present a prompt ablation in Appx. F.

**Grad-CAM [66] and Chefer *et al.* [12] overview.**

We begin by reviewing the Grad-CAM formulation, from which many gradient-based interpretability methods draw inspiration. Given a target class $c$, an input image $x$ and a model $f_\theta$, Grad-CAM 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 $\alpha_k^c$ are obtained from average pooling the gradients of the activations $A^k$:

$$\alpha_k^c = \mathsf{AveragePool}_{i,j}\left( \frac{\partial y^c}{\partial A_{ij}^k} \right). \quad (2)$$

The relevance map is then given by a linear combination of the activations, weighted by the importance from Eq. 2:

$$L_c = \mathsf{ReLU}\left( \sum_k \alpha_k^c A^k \right), \quad (3)$$

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

For ViTs [22], we follow the method of Chefer *et al.* [12]. Specifically, given the attention maps $A^k$ for each transformer block $k$, and relevance scores $R^k$ [48], the relevance map $L_c^{\mathsf{ViT}}$ is given by:

$$L_c^{\mathsf{ViT}} = \prod_k \bar{A}^k, \quad (4)$$

---

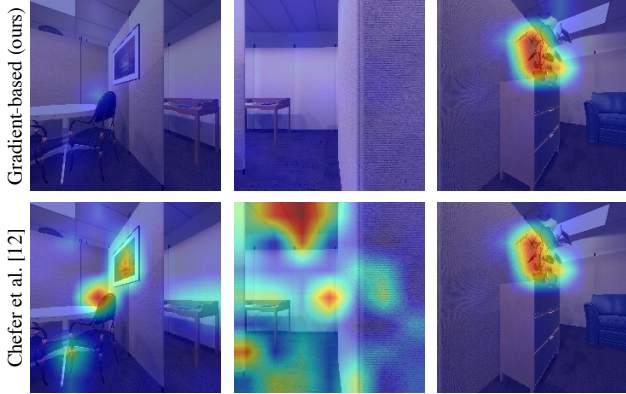[3]notebook exploring the effects of the 80 prompts on ImageNet-1k

Figure 7. **Gradient-based relevance visualization.** The target object is a plant. For our gradient-based strategy, derived from that of Chefer *et al.* [12], relevance is low when the object is not in the image and high otherwise. In contrast the original method, produces spurious relevance when the target object is not in the frame due to the normalization it employs. Notice when the plant is in the frame, the relevance map is similar.

$$\bar{A}^k = I + \mathbb{E}_h \left[ \mathrm{ReLU} \left( \nabla A^k \odot R^k \right) \right], \qquad (5)$$

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, $k = 1$.

It is standard to use an interpretability method like those of Selvaraju *et al.* [66] or Chefer *et al.* [12] to query for a class that is known to be in the image a priori. Hence, it is common practice to normalize relevance maps by subtracting the minimum relevance and normalizing by the difference between the maximum and minimum relevance. This results in a max value of 1.0. As mentioned in Sec. 4.3 this is not a suitable strategy for object navigation because in many frames the object is *not* in the image. Hence, in early experiments, we removed the normalization. Fig. 7 makes the differences between these strategies apparent. We notice this simple modification gives signal not only for true positive detections, but—critically—also for true negatives. Notice that when the plant is not in view, relevance is qualitatively low. However, when the plant is in view, relevance spikes in the region of the plant.

**MDETR overview.** MDETR [35] utilizes an encoder-decoder scheme to associate tokens in an input prompt referring to parts of an input image, with output boxes. MDETR utilizes a vision tower and a pre-trained language tower to project an image, text pair into a joint embedding space (similar to CLIP as discussed above). Image and text features are concatenated and passed to a transformer decoder head, which outputs boxes.The model is trained on a dataset of 200,000 images annotated with captions, boxes, and correspondence between words and boxes.

The model employs three losses during training. (1)

DETR bipartite matching loss without class information: Following DETR [7], each ground truth box is matched with its most similar predicted box and an L1 loss is applied. Unlike DETR, MDETR does not use any class information for the matching step. (2) A soft-token classification loss: when classifying a box, the target is a uniform distribution over all the tokens the box refers to and zero for other tokens. In this way the box is assigned to potentially many tokens depending on the ground truth. (3) A contrastive loss in the bottleneck embedding space: this is analogous to the CLIP loss discussed above.

To fine-tune MDETR for segmentation, the original box model is fine-tuned in two stages on the PhraseCut dataset [75]. First the model is fine-tuned for box prediction discussed above on PhraseCut data, which contains referring expressions and corresponding regions in the image (masks and boxes). The weights are frozen and a new segmentation head is trained using Dice/F1 loss and focal loss.

**OWL ViT overview.** OWL ViT [47] employs a two-stage training procedure to create an open-vocabulary model. During the first stage they train a CLIP-like model. However, while the original CLIP model uses the feature corresponding to a ViT [CLS] token to construct its multi-modal embedding space, OWL ViT instead pools over patch tokens to obtain an image representation. Intuitively, this encourages the image global feature to encode more local information from each patch. During the second stage of training, the pooling layer is removed. The patch tokens are passed to a linear projection head where they are then dotted against text features to determine class probabilities. An MLP box-projection head is introduced that predicts a box for each projected patch token. Note this process and losses for box fine-tuning are similar to that used by DETR [7]. Stage 1 is trained using 3.6 billion image-text pairs from the dataset used in LiT [80], using a standard CLIP loss. Stage 2 is fine-tuned using an agglomeration of existing box datasets with ~2M images total.

**Localization thresholds.** Each object localization method discussed in Sec. 4.3, requires a confidence threshold, which is standard when using a detector. To tune this threshold, we render 500 images in ROBOTHOR and 500 images in HABITAT MP3D with box annotations. Critically, we use the training rooms for these datasets, so none of the scenes overlap with those seen during downstream navigation testing. See Tab. 5 for thresholds, which are chosen to maximize an F1-score. Because PASTURE is a test set, we do *no* hyper-tuning on PASTURE. Instead we use the hyper-parameters from ROBOTHOR for PASTURE evaluations.

Here we present more information on computing our F1-scores. 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 $M_{\mathrm{pred}}^{o^+}$ a true positive (TP) if

| IDs | Localizer | Arch. | HABITAT | ROBOTHOR and PASTURE |
|---|---|---|---|---|
| ▲, △ | CLIP-Ref. | B/32 | – | 0.25 |
| ■, □ | CLIP-Ref. | B/16 | – | 0.125 |
| ▲, △ | CLIP-Patch | B/32 | – | 0.875 |
| ■, □ | CLIP-Patch | B/16 | – | 0.75 |
| ▲, △ | CLIP-Grad. | B/32 | 0.375 | 0.625 |
| ■, □ | CLIP-Grad. | B/16 | – | 0.375 |
| ◆, ◇ | MDETR | B3 | – | 0.95 |
| ▲, △ | OWL | B/32 | 0.2 | 0.125 |
| ■, □ | OWL | B/16 | – | 0.125 |

Table 5. **Object localization hyperparameters.** Hyperparameters returned from a grid-search on object localization performance on HABITAT and ROBOTHOR train scenes. Note, no hyperparameter tuning was conducted on PASTURE, which is strictly a test set. Missing entries indicate that we did not evaluate these models due to lacking performance on other datasets (i.e., on ROBOTHOR).

$\mathsf{score}^+ = \sum(M^{o^+}_{\text{pred}} \odot M^{o^+}_{\text{GT}})/\sum M^{o^+}_{\text{pred}} > 0.5$, where $o^+ \in \mathcal{O}^+$, $M^{o^+}_{\text{GT}}$ 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 $0 < \mathsf{score}^+ \leq 0.5$ or $\sum M^{o^-}_{\text{pred}} > 0$, where $o^- \in \mathcal{O}^-$. False negatives (FN) arise when $\sum M^{o^+} = 0$. F1 is computed in the standard way: $\text{TP}/(\text{TP} + 0.5(\text{FP} + \text{FN}))$. 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.

## D. Additional Visualization

To give a more qualitative perspective on our results, we additionally visualize some key aspects of the method and also evaluation trajectories. In Fig. 8, we show back-projecting 2D object relevancy to 3D, which is a key part of the CoW pipeline. We provide success and failure sample trajectories in Fig. 9.

## E. Dataset Details

**ROBOTHOR and HABITAT MP3D dataset details.** There are 11 HABITAT MP3D and 15 ROBOTHOR test scenes, in the official validation sets, which we use as our test sets. HABITAT includes 2,195 evaluation episodes, while ROBOTHOR has 1,800. We consider the following 21 categories in HABITAT: `chair`, `table`, `picture`, `cabinet`, `cushion`, `sofa`, `bed`,

(a) Egocentric RGB  (b) CLIP-Grad. B/32 relevance overlay  (c) Depth projected relevance (aggregated)
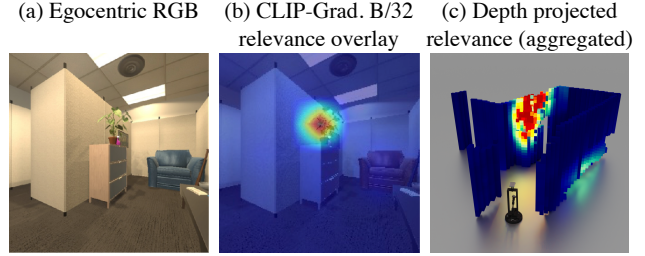


Figure 8. **Projection of object relevance.** (a) Egocentric RGB. Note, a CoW also receives a depth image. (b) Raw CLIP-Grad. B/32 prediction for the image targeting the "plant" class. (c) Back-projection of object relevance, aggregated over time, into a 3D map using agent pose estimates. Areas of high relevence make natural targets for navigation.

`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 MP3D. We provide testing episode counts per category in Fig. 10.

**PASTURE additional statistics.** In Fig. 11, we show word counts for various color, material, and spatial attributes for the appearance and spatial class remapping. In Fig. 12, we repeat similar analysis for hidden objects, reporting the word frequency for objects that contain the target objects.

**PASTURE uncommon object asset licensing.** We found the base instances for the PASTURE uncommon object split from CGTrader, an online repository where CAD hobbyists, artists, and professionals host their models. We choosing our objects, we selected 12 instances that had a "Royalty Free License" and were free to download at the time of dataset construction. We acknowledge all artists and provide links to their work:

- "tie-dye surfboard" by *adhil3dartist* and re-textured to be rainbow
- "whiteboard saying CVPR" by *w-stone-art* and re-textured to say CVPR
- "llama wicker basket" by *eelh*
- "green plastic crate" by *Snowdrop-2018* and modified to include only the green create
- "rice cooker" by *fleigh* and re-textured
- "mate gourd" by *mcgamescompany*
- "red and blue tricycle" by *POLY1PROPS*

Figure 9. **Trajectory visualization.** Frames are egocentric views. Color indicates trajectory progress, where blue indicating trajectory start and white indicating trajectory end. Target objects are boxed in green, while distractor objects are boxed in red.



Figure 10. HABITAT **MP3D and** ROBOTHOR **episode splits.** Distribution of episodes in each CVPR 2021 object navigation challenge validation set that we adopt as our test set.

- "white electric guitar" by *demolitions2000*
- "espresso machine" by *WolfgangNikolas*
- "wooden toy airplane" by *fomenos*
- "gingerbread house" by *Empire-Assets*
- "graphics card" by *Biggie-3D*

PASTURE **sample prompts.** Here we provide some sample prompts for our appearance, spatial, and hidden object instance remapping.

- appearance remapping: from "spray bottle" to "small, green, plastic spray bottle"

- spatial remapping: from "spray bottle" to "spray bottle on a coffee table near a house plant"

- hidden remapping: from "spray bottle" to "spray bottle under the bed"

**Retrieval of uncommon objects on LAION-5B [64]**
We include sample data from CLIP retrieval for both ROBOTHOR and PASTURE uncommon objects in Fig. 13. We included the results to show that CLIP is familiar with the concepts that we chose when creating the uncommon split of PASTURE and that qualitative image results rival those of the more common ROBOTHOR objects.

## F. Prompt Ensemble Ablation

To verify the benefits of using the 80 prompt ensemble created by OpenAI—denoted as *prompt ens.*—, we compare performance for two models on ROBOTHOR. The competing approach uses a single prompt for each class

Figure 11. **Attribute distribution for PASTURE.** (a) Distribution of color and material attribute word frequency in the PASTURE appearance captions. (b) Distribution of on/in and near attributes by frequency in the PASTURE spatial captions.



Figure 12. **Hidden reference object distribution for PASTURE.** Word frequency of large objects that target objects are hidden "in" or "under" in the PASTURE hidden object captions. Names are based on minimal descriptions needed to identify the object in the room. For example, "brown sofa" vs. "white sofa".

"a photo of a {}."—denoted as *photo*. As we see in Tab. 6, prompt ens. boosts SPL by by 0.2 for the CLIP-Grad. model. These, results suggest that the two variations of prompting strategies have marginal influence on down-

| | | CoW breeds | | | ROBOTHOR | |
| ID | Loc. | Arch. | Post | Exp. Strategy | SPL | SR |
| △ | CLIP-Grad. | B/32 | ✓ | photo | 9.5 | **15.2** |
| △ | CLIP-Grad. | B/32 | ✓ | prompt ens. | **9.7** | **15.2** |

Table 6. **Prompt ablation.** For a fixed object localizers (CLIP-Grad. B/32 with post processing), we ablate over different choices of prompts. We find that the 80 prompt ensemble (prompt ens.), introduced by OpenAI, outperforms the simple prompt: "a photo of a {}." (photo) in most cases. However, deltas are not large, suggesting that this is a less critical design decision in the CoW framework.

stream performance. For all experiments in the main paper that use CLIP, we use the 80 prompt ensemble in conjunction with the class label specific for a task (e.g., "spray bottle under the bed").

## G. Category-level Results

For completeness we also include salient category-level results for an OWL (△) and CLIP-Grad. (△) B/32 models with post-processing for HABITAT MP3D, ROBOTHOR, and PASTURE. For a comparison for the appearance (Appear.) and spatial (Space) PASTURE splits, see Tab. 7. For Appear. and Space with distractors (distract), see Tab. 8.

## RoboTHOR Objects



| "garbage can" | "laptop" | "baseball bat" | "mug" | "bowl" | "alarm clock" |
| "apple" | "basketball" | "TV" | "vase" | "plant" | "spray bottle" |

## Pasture Uncommon Objects

| "llama wicker basket" | "red and blue tricycle" | "tie-dye surfboard" | "wooden toy airplane" | "green plastic crate" | "white electric guitar" |
| "whiteboard saying CVPR" | "graphics card" | "gingerbread house" | "espresso machine" | "rice cooker" | "mate gourd" |

Figure 13. **Qualitative CLIP retrieval.** We conduct CLIP retrieval on LAION-5B [64] using the class names shown in quotes with the interface provided here. We show results for both ROBOTHOR and PASTURE uncommon objects. While CLIP is not trained on LAION-5B, this figure is included to give an idea of the type of noisy internet image-text data is trained on. Retrieval returns reasonable results for uncommon objects, suggesting CLIP is able to semantically distinguish these objects.

For hidden object with and without distractors see Tab. 9. For uncommon objects see Tab. 10. For ROBOTHOR see Tab. 11. For HABITAT see Tab. 12.

## H. Additional Failure Analysis

For addition failure analysis on CLIP-Ref., CLIP-Patch, and CLIP-Grad. see Fig. 14. All models have a ViT-B/32 architecture and apply post-processing. We notice that that CLIP-Patch and CLIP-Grad. have a higher fraction of object localization failure when compared to OWL. For CLIP-Ref., where performance is in the very low success regime (e.g., $< 3\%$), we notice less consistent patterns.
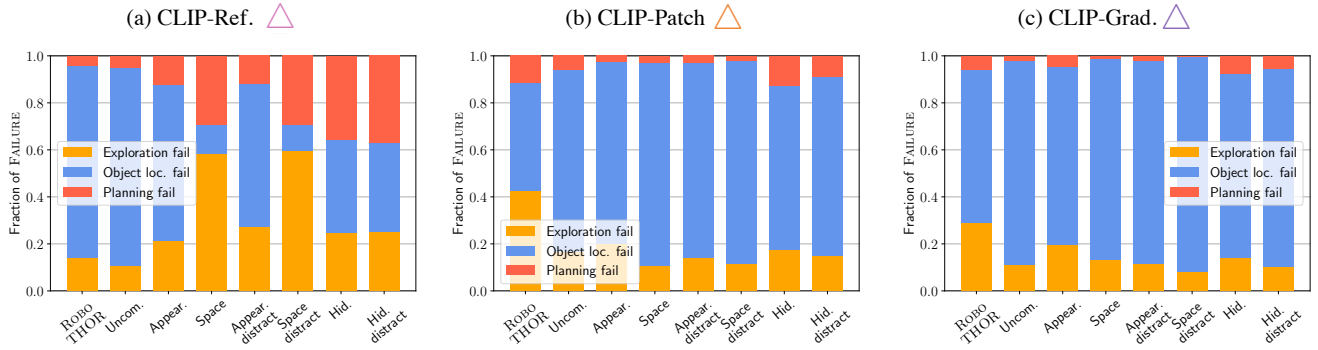
Figure 14. **Failure analysis for more models.** (a, b, c) All show additional error analysis. When comparing CLIP-Patch and CLIP-Grad. to OWL shown in Fig. 6, we notice that the former have a higher percentage of object localization failures and also lower success rate downstream on average.

| category | PASTURE Appear. | | | | PASTURE Space | | | |
|---|---|---|---|---|---|---|---|---|
| | SR | SPL | SR | SPL | SR | SPL | SR | SPL |
| ALARMCLOCK | 6.7 | 4.0 | **23.3** | **10.7** | 3.3 | 3.3 | **10.0** | **3.8** |
| APPLE | 6.7 | 5.7 | **36.7** | **17.0** | **10.0** | **8.4** | 3.3 | 1.0 |
| BASEBALLBAT | 0.0 | 0.0 | **3.3** | **1.2** | 3.3 | 2.5 | **6.7** | **2.7** |
| BASKETBALL | 6.7 | 2.8 | **36.7** | **24.1** | 10.0 | 5.6 | **36.7** | **26.8** |
| BOWL | 3.3 | 0.5 | **13.3** | **5.9** | 10.0 | 5.6 | **16.7** | **6.9** |
| GARBAGECAN | 26.7 | 20.2 | **50.0** | **31.5** | 30.0 | 23.0 | **40.0** | **23.2** |
| HOUSEPLANT | 20.0 | 16.9 | **30.0** | **20.2** | 13.3 | 10.8 | **40.0** | **21.9** |
| LAPTOP | 13.3 | 10.6 | **20.0** | **11.5** | 13.3 | 9.6 | **20.0** | **13.7** |
| MUG | 10.0 | 7.5 | **46.7** | **27.4** | 10.0 | **7.5** | **13.3** | 5.4 |
| SPRAYBOTTLE | 16.7 | 13.6 | **33.3** | **19.2** | **16.7** | **15.8** | 16.7 | 6.8 |
| TELEVISION | 10.0 | **10.0** | **13.3** | 8.9 | 6.7 | 6.4 | **20.0** | **9.9** |
| VASE | **23.3** | **17.5** | 10.0 | 9.1 | **13.3** | **9.8** | 10.0 | 5.4 |

Table 7. **Attribute object navigation.** Appearance-based captions consistently perform better than spatial captions. OWL consistently performs better than CLIP-Grad.

| category | PASTURE Appear. distract | | | | PASTURE Space distract | | | |
|---|---|---|---|---|---|---|---|---|
| | SR | SPL | SR | SPL | SR | SPL | SR | SPL |
| ALARMCLOCK | 3.3 | 3.0 | **13.3** | **6.5** | **6.7** | **6.3** | 6.7 | 2.8 |
| APPLE | **10.0** | 6.4 | 10.0 | **7.4** | 3.3 | 3.3 | **10.0** | **4.0** |
| BASEBALLBAT | 0.0 | 0.0 | **13.3** | **4.5** | 0.0 | 0.0 | **10.0** | **8.1** |
| BASKETBALL | 6.7 | 3.3 | **20.0** | **12.6** | 16.7 | 9.4 | **16.7** | **9.7** |
| BOWL | 3.3 | 3.2 | **16.7** | **8.5** | 10.0 | 8.6 | **23.3** | **12.6** |
| GARBAGECAN | 26.7 | 19.9 | **30.0** | **21.6** | 13.3 | 10.5 | **26.7** | **18.2** |
| HOUSEPLANT | 10.0 | 6.0 | **16.7** | **11.0** | 13.3 | 10.8 | **23.3** | **13.7** |
| LAPTOP | 16.7 | **13.6** | **23.3** | 11.9 | **16.7** | **11.9** | 16.7 | 11.8 |
| MUG | 6.7 | 5.1 | **26.7** | **17.8** | **10.0** | **7.8** | 6.7 | 2.7 |
| SPRAYBOTTLE | 13.3 | 12.4 | **26.7** | **15.2** | 16.7 | **15.4** | **20.0** | 8.0 |
| TELEVISION | 6.7 | 6.6 | **26.7** | **15.5** | 6.7 | 3.3 | **20.0** | **12.8** |
| VASE | **13.3** | **10.2** | 10.0 | 8.6 | 10.0 | 6.5 | **13.3** | **8.4** |

Table 8. **Attribute object navigation with distractors.** Distractors consistently hurt performance compared to the no distractor numbers in Tab. 7, suggesting that models cannot make full use of remapped classes with attributes.

| category | PASTURE Hidden | | | | PASTURE Hidden distract | | | |
|---|---|---|---|---|---|---|---|---|
| | SR | SPL | SR | SPL | SR | SPL | SR | SPL |
| AlarmClock | 26.7 | **15.8** | 6.7 | 3.0 | 20.0 | **9.9** | 23.3 | 8.3 |
| Apple | 10.0 | **9.5** | 13.3 | 8.9 | 6.7 | 4.7 | **10.0** | 7.9 |
| BaseballBat | 13.3 | **9.2** | 13.3 | 5.8 | 10.0 | 7.2 | 3.3 | 2.9 |
| BasketBall | 10.0 | 5.1 | 26.7 | 15.2 | 16.7 | 12.2 | **20.0** | **15.0** |
| Bowl | 16.7 | 9.2 | 23.3 | 11.9 | 23.3 | 15.1 | 20.0 | 10.6 |
| GarbageCan | 13.3 | 7.8 | 26.7 | 13.8 | 13.3 | 7.0 | 16.7 | 10.9 |
| HousePlant | 13.3 | 9.6 | 16.7 | 10.8 | 16.7 | 10.9 | 16.7 | 11.6 |
| Laptop | **10.0** | 8.1 | 10.0 | 8.4 | 26.7 | 20.7 | 10.0 | 8.4 |
| Mug | 13.3 | 7.3 | 33.3 | 23.8 | 20.0 | 12.5 | 23.3 | 18.4 |
| SprayBottle | **13.3** | **8.4** | 0.0 | 0.0 | **20.0** | **10.3** | 0.0 | 0.0 |
| Television | 16.7 | 8.0 | 36.7 | 22.9 | 10.0 | 6.9 | 16.7 | 11.4 |
| Vase | 16.7 | **11.9** | 23.3 | 11.2 | 10.0 | 7.0 | 13.3 | **7.0** |

Table 9. **Hidden object navigation category-level results.**

| category | PASTURE Uncom. | | | |
|---|---|---|---|---|
| | SR | SPL | SR | SPL |
| GingerbreadHouse | 20.0 | 14.3 | 26.7 | 18.6 |
| EspressoMachine | 10.0 | 7.7 | 46.7 | 24.6 |
| Crate | 23.3 | 18.2 | 40.0 | 27.0 |
| ElectricGuitar | 16.7 | 10.0 | 46.7 | 30.8 |
| RiceCooker | 3.3 | 2.9 | 20.0 | 11.6 |
| LlamaWickerBasket | 16.7 | 12.6 | 30.0 | 24.5 |
| Whiteboard | **63.3** | **43.2** | 30.0 | 18.7 |
| Surfboard | 26.7 | 20.6 | 60.0 | 38.9 |
| Tricycle | 10.0 | 9.0 | 53.3 | 31.7 |
| GraphicsCard | 3.3 | 2.1 | 13.3 | 6.0 |
| Mate | **0.0** | **0.0** | **0.0** | **0.0** |
| ToyAirplane | 0.0 | 0.0 | 26.7 | 13.7 |

Table 10. **Uncommon object navigation category-level results.**

| category | ROBOTHOR | | | |
|---|---|---|---|---|
| | SR | SPL | SR | SPL |
| AlarmClock | 5.3 | 3.1 | **30.7** | **19.2** |
| Apple | 15.3 | 9.7 | **34.0** | **19.6** |
| BaseballBat | **4.0** | **1.5** | 2.0 | 0.5 |
| BasketBall | 19.3 | 14.8 | **36.0** | **25.3** |
| Bowl | 5.3 | 4.0 | **18.0** | **11.1** |
| GarbageCan | 30.0 | 21.0 | **50.0** | **32.2** |
| HousePlant | 30.7 | 18.7 | **36.7** | **24.8** |
| Laptop | 16.0 | 10.6 | **20.0** | **11.3** |
| Mug | 10.7 | 5.0 | **40.7** | **25.5** |
| SprayBottle | 8.0 | 5.2 | **23.3** | **13.8** |
| Television | **29.3** | **16.9** | 23.3 | 13.4 |
| Vase | **8.0** | **5.8** | 6.0 | 5.7 |

Table 11. **ROBOTHOR category-level results.**

| category | PASTURE HABITAT | | | |
|---|---|---|---|---|
| | SR | SPL | SR | SPL |
| CHAIR | 5.1 | 2.2 | **7.4** | **3.7** |
| TABLE | **34.7** | **18.4** | 21.2 | 9.6 |
| PICTURE | 1.2 | 0.5 | **1.9** | **1.1** |
| CABINET | **9.4** | 4.7 | 8.3 | **4.9** |
| CUSHION | 5.0 | 3.0 | **12.1** | **5.8** |
| SOFA | **0.0** | **0.0** | **0.0** | **0.0** |
| BED | **0.0** | **0.0** | **0.0** | **0.0** |
| CHEST_OF_DRAWERS | 0.0 | 0.0 | **1.6** | **0.4** |
| PLANT | **5.7** | **3.6** | 2.3 | 1.3 |
| SINK | 2.4 | 1.6 | **2.4** | **1.8** |
| TOILET | **0.0** | **0.0** | **0.0** | **0.0** |
| STOOL | 0.0 | 0.0 | **2.9** | **2.4** |
| TOWEL | 0.0 | 0.0 | **1.4** | **0.7** |
| TV_MONITOR | **0.0** | **0.0** | **0.0** | **0.0** |
| SHOWER | **7.1** | **4.4** | 1.4 | 0.9 |
| BATHTUB | 0.0 | 0.0 | **11.1** | **3.1** |
| COUNTER | **6.1** | **3.1** | 0.0 | 0.0 |
| FIREPLACE | **10.0** | **5.0** | 3.3 | 1.7 |
| GYM_EQUIPMENT | **0.0** | **0.0** | **0.0** | **0.0** |
| SEATING | **12.6** | **5.2** | 0.0 | 0.0 |
| CLOTHES | **8.0** | **4.3** | 4.0 | 2.1 |

Table 12. **HABITAT category-level results.**