
A Simple Approach for Visual Rearrangement: 3D Mapping and Semantic Search

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Physically rearranging objects is an important capability for embodied agents.
2 Visual room rearrangement evaluates an agent’s ability to rearrange objects in a
3 room to a desired goal based solely on visual input. We propose a simple yet
4 effective method for this problem: (1) search for and map which objects need
5 to be rearranged, and (2) rearrange each object until the task is complete. Our
6 approach consists of an off-the-shelf semantic segmentation model, voxel-based
7 semantic map, and semantic search policy to efficiently find objects that need to be
8 rearranged. On the AI2-THOR Rearrangement Challenge, our method improves
9 on current state-of-the-art end-to-end reinforcement learning-based methods that
10 learn visual rearrangement policies from 0.53% correct rearrangement to 15.11%,
11 using only 2.7% as many samples from the environment.

12 1 Introduction

13 Physically rearranging objects is an everyday skill for humans, but remains a core challenge for
14 embodied agents that assist humans in realistic environments. Natural environments for humans
15 are complex and require generalization to a combinatorially large number of object configurations
16 [Batra et al., 2020a]. Generalization in complex realistic environments remains an immense practical
17 challenge for embodied agents, and the rearrangement setting provides a rich test bed for embodied
18 generalization in these environments. The rearrangement setting combines two challenging perception
19 and control tasks: (1) understanding the state of a dynamic 3D environment, and (2) acting over a
20 long horizon to reach a goal. These problems have traditionally been studied independently by the
21 vision and reinforcement learning communities [Chaplot et al., 2021], but the advent of large models
22 and challenging benchmarks is showing that both components are important for embodied agents.

23 Reinforcement learning (RL) can excel at embodied tasks, especially if centuries of experience
24 can be leveraged [Weihs et al., 2021, Chaplot et al., 2020b, Ye et al., 2021] for training. In a
25 simulated environment with unlimited retries, this experience is cheap to obtain, and agents can
26 explore randomly until a good solution is discovered by the agent. This pipeline works incredibly well
27 for tasks like point navigation [Wijmans et al., 2020], but in some cases this strategy is not enough.
28 As the difficulty of embodied learning tasks increases, the agent must generalize to an increasing
29 number of environment configurations, and broadly scaled experience can become insufficient.

30 In the rearrangement setting, a perfect understanding of the environment simplifies the problem: an
31 object is here, it should go there, and the rest can be solved with grasping and planning routines.
32 Representing the information about the locations and states of objects in an accessible format is
33 therefore an important contribution for the rearrangement setting. Our initial experiments suggest
34 that accurate 3D semantic maps of the environment are one such accessible format for visual
35 rearrangement. With accurate 3D semantic maps, our method rearranges 15.11% of objects correctly,
36 and requires significantly less experience from the environment to do so. While end-to-end RL

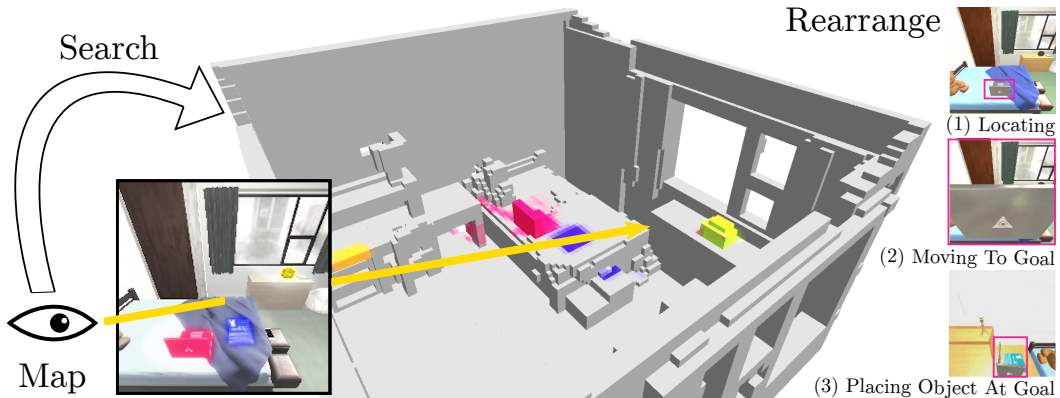


Figure 1: Our method incrementally builds voxel-based *Semantic Maps* from visual observations and efficiently finds objects using a *Semantic Search Policy*. We visualize an example rearrangement on the right with the initial position of the pink object (laptop on the bed), followed by the agent holding the object (laptop), and finally the destination position of the object (laptop on the desk).

37 requires up to 75 million environment steps in Weihs et al. [2021], our method only requires 2 million
 38 samples and trains offline. Our results suggest end-to-end RL without an accurate representation of
 39 the scene may be missing out on a fundamental aspect of understanding of the environment.

40 We demonstrate how semantic maps help agents effectively understand dynamic 3D environments
 41 and perform visual rearrangement. These dynamic environments have elements that can move (like
 42 furniture), and objects with changing states (like the door of a cabinet). We present a method that
 43 builds accurate semantic maps in these dynamic environments, and reasons about what has changed.
 44 Deviating from prior work that leverages end-to-end RL, we propose a simple approach for visual
 45 rearrangement: (1) search for and map which objects need to be rearranged, and (2) procedurally
 46 rearrange objects until a desired goal configuration is reached. We evaluate our approach on the
 47 AI2-THOR Rearrangement Challenge [Weihs et al., 2021] and establish a new state-of-the-art.

48 We propose an architecture for visual rearrangement that builds voxel-based semantic maps of the
 49 environment and rapidly finds objects using a search-based policy. Our method shows an improvement
 50 of 14.72 absolute percentage points over current work in visual rearrangement, and is robust to the
 51 accuracy of the perception model, the budget for exploration, and the size of objects being rearranged.
 52 We conduct ablations to diagnose where the bottlenecks are for visual rearrangement, and find that
 53 accurate scene understanding is the most crucial. As an upper bound, when provided with a perfect
 54 semantic map, our method solves 38.33% of tasks, a potential for significant *out-of-the-box* gains as
 55 better perception models are developed. Our results show the importance of building effective scene
 56 representations for embodied agents in complex and dynamic visual environments.

57 2 Related Work

58 **Embodied 3D Scene Understanding.** Knowledge of the 3D environment is at the heart of various
 59 tasks for embodied agents, such as point navigation [Anderson et al., 2018a], image navigation [Batra
 60 et al., 2020b, Yang et al., 2019], vision language navigation [Anderson et al., 2018b, Shridhar et al.,
 61 2020], embodied question answering [Gordon et al., 2018, Das et al., 2018], and more. These tasks
 62 require an agent to reason about its 3D environment. For example, vision language navigation [An-
 63 derson et al., 2018b, Shridhar et al., 2020] requires grounding language in an environment goal,
 64 and reasoning about where to navigate and what to modify in the environment to reach that goal.
 65 Reasoning about the 3D environment is especially important for the rearrangement setting, and has a
 66 rich interdisciplinary history in the robotics, vision, and reinforcement learning communities.

67 **Visual Room Rearrangement.** Rearrangement has long been one of the fundamental tasks in
 68 robotics research [Ben-Shahar and Rivlin, 1996, Stilman et al., 2007, King et al., 2016, Krontiris and
 69 Bekris, 2016, Yuan et al., 2018, Correll et al., 2018, Labbé et al., 2020]. Typically, these methods
 70 address the challenge in the context of the state of the objects being fully observed [Cosgun et al.,
 71 2011, King et al., 2016], which allows for efficient and accurate planning-based solutions. In contrast,
 72 there has been recent interest in room rearrangement inside a realistic 3D simulator [Batra et al.,

73 2020a, Weihs et al., 2021, Gadre et al., 2022] where the states of objects and the rearrangement goal
74 are not directly observed. In these cases, the simulator only provides a direct visual input, and the
75 simulated environment is relatively complex and realistic. This latest iteration of rearrangement shares
76 similarity with various other challenging embodied AI tasks such as embodied navigation [Anderson
77 et al., 2018a, Batra et al., 2020b, Chaplot et al., 2020a, Shridhar et al., 2020, Francis et al., 2021, Min
78 et al., 2021, Pashevich et al., 2021, Singh et al., 2021] and embodied question answering [Gordon
79 et al., 2018, Das et al., 2018], which require finding objects and reasoning about their state.

80 **AI2-THOR Rearrangement Challenge.** Our work builds on the latest rearrangement methods
81 and demonstrates how building accurate voxel-based semantic maps can produce significant gains.
82 We focus on the AI2-THOR Rearrangement Challenge [Weihs et al., 2021], which uses AI2-THOR,
83 an open-source and high-fidelity simulator used in many prior works [Gadre et al., 2022, Weihs
84 et al., 2021, Shridhar et al., 2020, Gordon et al., 2018]. Prior works on this challenge have studied a
85 variety of approaches, including end-to-end RL in Weihs et al. [2021], and a planning-based approach
86 in Gadre et al. [2022]. Our approach is the first to use voxel-based semantic maps to infer what
87 to rearrange from an experience goal as described by Batra et al. [2020a]. Though both Gadre
88 et al. [2022] and our method use planning, Gadre et al. [2022] use a graph-based continuous scene
89 representation, and we use voxel-based semantic maps instead, which we show is more effective.

90 **3D Mapping & Search.** Agents that interact with an embodied world through navigation and ma-
91 nipulation must keep track of the world (mapping) [Thrun, 2002] and themselves (localization) [Thrun
92 et al., 2001]—both extensively studied in robotics by processing low-level information [Engel et al.,
93 2014], building semantic maps [Kuipers and Byun, 1991] and more recently, via techniques specifi-
94 cally developed to handle dynamic and general aspects of the environment [Rünz and Agapito, 2017,
95 Rosinol et al., 2021, Wong et al., 2021]. When semantics are more important than precision, such
96 as for embodied learning tasks, recent methods have looked at neural network-based maps [Gupta
97 et al., 2017, Chen et al., 2019, Wu et al., 2019b, Chaplot et al., 2020b, Blukis et al., 2021, Chaplot
98 et al., 2021]. Our method builds on these and adopts the use of a voxel-based semantic map and pre-
99 trained semantic segmentation model—a similar methodological setup to Chaplot et al. [2021], Min
100 et al. [2021]. However, our method diverges from these prior works by using multiple voxel-based
101 semantic maps to infer what to rearrange from an experience goal as described by Batra et al. [2020a].
102 These prior works have instead considered geometric goals in Chaplot et al. [2021] and language
103 goals in Min et al. [2021], and ours is the first to consider an experience goal [Batra et al., 2020a].
104 Furthermore, while a search-based policy is used in Min et al. [2021], we are the first to use search
105 with an unspecified destination (ie, the agent does not know what kind of object is it looking for).

106 3 Methodology

107 In this section, we present a simple approach for solving visual rearrangement problems. We begin the
108 section by discussing the visual rearrangement problem statement and metrics we use for evaluation.
109 We then discuss our methodological contributions. First, we propose to build multiple voxel-based
110 semantic maps representing the environment in different configurations. Second, we propose a policy
111 that efficiently finds objects that need to be rearranged. Third, we propose a method for inferring the
112 rearrangement goal from two semantic maps to efficiently solve visual rearrangement tasks.

113 **Visual rearrangement definition and evaluation metrics.** Consider the rearrangement setting
114 defined by Batra et al. [2020a], which is a special case of a Markov Decision Process (MDP)
115 augmented with a goal specification $g = \phi(s_0, S^*)$. This goal specification encodes the set of states
116 S^* for which the rearrangement task is considered solved from initial state s_0 . The agent typically
117 does not directly observe the set of goal states S^* , and this is reflected by the goal specification
118 function $\phi : S \times 2^S \rightarrow \mathcal{G}$. We consider a setting where the rearrangement goal g is specified
119 visually and the agent initially observes the environment in its goal configuration. This setting is
120 especially challenging because the agent must remember what the environment initially looked like to
121 infer the set of goal states. Once the goal has been understood and rearrangement has been attempted,
122 we evaluate agents using metrics introduced by Weihs et al. [2021]. We consider a *Success* metric
123 that measures the proportion of tasks for which the agent has correctly rearranged all objects and
124 misplaced none during rearrangement. This metric is strict in the sense that an agent receives a
125 success of 0.0 if at least one object is misplaced—even if all others are correctly rearranged. We

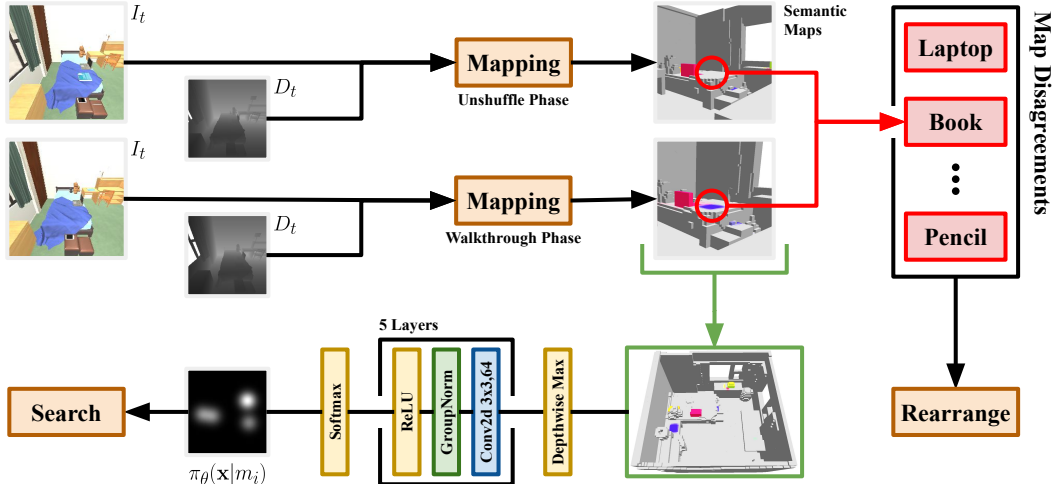


Figure 2: Overview of our method for an example task. Our method incrementally builds voxel-based *Semantic Maps* from visual observations. Our *Semantic Search Policy* helps build accurate maps by selecting navigation goals to efficiently find objects that need to be rearranged. Once accurate maps are built, our method compares the *Semantic Maps* to identify disagreements between the maps, and rearranges objects to resolve those disagreements using a deterministic rearrangement policy.

126 consider an additional *%Fixed Strict* metric that measures the proportion of objects per task correctly
 127 rearranged, equal to 0.0 per task if any were misplaced. This second metric is more informative
 128 regarding how close the agent was to solving each task. Effective agents will correctly rearrange all
 129 objects in the scene to their goal configurations, maximizing their *Success* and *%Fixed Strict*.

130 **Building two semantic maps.** Our approach builds off recent work that uses voxel-based semantic
 131 maps in embodied settings [Min et al., 2021, Chaplot et al., 2021]. Our work differs from these in
 132 that we use multiple voxel-based semantic maps to encode both the goal state and current state of the
 133 environment. In particular, we build two semantic maps $m_0, m_1 \in \mathcal{R}^{H \times W \times D \times C}$ that represent 3D
 134 grids with $H \times W \times D$ voxels. Each voxel is represented with a categorical distribution on C classes
 135 encoding which class is likely to occupy each voxel. Empty voxels are assigned the zero vector. In an
 136 initial observation phase for each task, our agent navigates the scene and builds m_0 , a semantic map
 137 encoding the goal configurations for objects in the scene. Likewise, in a second interaction phase, our
 138 agent navigates the scene and builds m_1 , a semantic map encoding the current state of objects in the
 139 scene. At every timestep during each phase, pose, RGB, and depth images are observed, and either
 140 m_0 or m_1 is updated depending on which instance of the scene the agent is currently observing.

141 **Incorporating semantic predictions in the maps.** Each semantic map is initialized to all zeros
 142 and, at every timestep t , semantic predictions from Mask R-CNN [He et al., 2017] are added to
 143 the map. Given the RGB image observation I_t , we generate semantic predictions from Mask R-
 144 CNN consisting of the probability of each pixel belonging to a particular class. We filter these
 145 predictions to remove those with a detection confidence lower than 0.9 and conduct an ablation in
 146 Section 4.3. We follow Chaplot et al. [2021] and generate an egocentric point cloud c_t^{ego} using the
 147 depth observation D_t . Each point in this point cloud is associated with a pixel in the image I_t and
 148 a vector of class probabilities from Mask R-CNN. Given the current pose x_t , we then transform
 149 the egocentric point cloud c_t^{ego} from the agent’s coordinate system to world coordinate system.
 150 This transformation results in a geocentric point cloud c_t^{geo} that is converted to a geocentric voxel
 151 representation $v_t^{geo} \in \mathcal{R}^{H \times W \times D \times C}$ of the same cardinality as the semantic maps. We additionally
 152 generate a voxelized mask $v_t^{mask} \in \mathcal{R}^{H \times W \times D \times 1}$ that equals one for every occupied voxel in v_t^{geo}
 153 and zero otherwise. New semantic predictions are added to the maps with a moving average.

$$m_i[t + 1] = m_i[t] \odot (1 - v_t^{mask}(1 - \epsilon)) + v_t^{geo}(1 - \epsilon) \quad (1)$$

154 The update in Equation 1 allows voxels to be updated at different rates depending on how frequently
 155 they are observed. The hyperparameter $\epsilon \in (0, 1)$ controls how quickly the semantic maps are
 156 updated to account for new semantic predictions, and is set to 0.5 in our experiments. An overview

Algorithm 1 3D Mapping and Semantic Search For Visual Rearrangement

Require: visual rearrangement environment e , initial voxel-based semantic maps $m_0, m_1 \in \mathcal{R}^{H \times W \times D \times C}$, search-based policy $\pi_\theta(\mathbf{x}|m)$, pre-trained semantic segmentation model g

for each phase $i \in \{0, 1\}$ **do**

for each I_t, D_t, x_t observed **do**

$v_t^{geo}, v_t^{mask} \leftarrow \text{project}(g(I_t), D_t, x_t)$ ▷ project to voxels

$m_i[t] \leftarrow m_i[t-1] \odot (1 - v_t^{mask}(1 - \epsilon)) + v_t^{geo}(1 - \epsilon)$ ▷ update map

if goal is reached or goal does not exist **then**

 goal $\sim \pi_\theta(\mathbf{x}|m_i[t])$ ▷ emit a semantic search goal

end if

 navigate to goal

end for

end for

while a disagreement d between m_0 and m_1 is detected **do**

 navigate to d in m_1 and rearrange d to match m_0

end while

157 of how our two semantic maps are built is shown in Figure 2. We’ve detailed how the semantic maps
158 are constructed from observations, and we will next describe how navigation goals are selected.

159 **Locating objects with a search-based policy.** Building accurate maps requires locating and
160 observing every object in the scene so they can be added to the maps. This requires intelligently
161 selecting navigation goals based on where objects are likely to be. We learn a high-level policy
162 $\pi_\theta(\mathbf{x}|m_i)$ that builds off recent work in Min et al. [2021], Chaplot et al. [2021] and parameterizes a
163 distribution over 3D search locations in the environment. The input to the policy is a 3D semantic
164 map m_i from whichever phase is currently active. The policy is a 5-layer 2D convolutional neural
165 network that processes a 3D semantic map m_i and outputs a categorical distribution over voxels in m_i ,
166 corresponding to 3D search locations. The policy is trained using maximum likelihood training with
167 an expert distribution $p^*(\mathbf{x})$ that captures the locations of the K objects the agent should rearrange in
168 the current scene. This expert distribution in Equation 2 is a Gaussian mixture model with a mode
169 centered at the location μ_k of each object, and a variance hyperparameter σ^2 for each mode.

$$p^*(\mathbf{x}) \propto \frac{1}{K} \sum_{k=1}^K \mathcal{N}(\mathbf{x}; \mu_k, \sigma^2 I) \quad (2)$$

170 Once a policy $\pi_\theta(\mathbf{x}|m_i)$ is trained that captures a semantic prior for object locations, we use planning
171 to reach goals sampled from the policy. We build a planar graph that represents traversable space
172 derived from voxel occupancy in the semantic map, and use Dijkstra’s algorithm [Dijkstra, 1959]
173 to find the shortest path from the agent’s current location to the goal. We filter navigation goals to
174 ensure only feasible goals are sampled, and then allow sufficient time for each navigation goal to be
175 reached. Once the current goal is reached, we sample another goal and call the planner again.

176 **Inferring the rearrangement goal from the maps.** Once two semantic maps are built, we compare
177 them to extract differences in object locations, which we refer to as map disagreements. These
178 disagreements represent objects that need to be rearranged by the agent. To locate disagreements,
179 we first use OpenCV [Bradski, 2000] to label connected voxels of the same class as object instances.
180 We consider voxels with nonzero probability of class c to contain an instance of that class. Object
181 instances are then matched between phases by taking the assignment of object instances that minimizes
182 the difference in appearance between instances of the same class. We leverage the Hungarian
183 algorithm [Kuhn and Yaw, 1955], and represent appearance by the average color of an object instance
184 in the map. Once objects are matched, we label pairs separated by > 0.05 meters as disagreements.
185 Given a set of map disagreements $\{(x_1, x_1^*), (x_2, x_2^*), \dots, (x_N, x_N^*)\}$ represented by the current pose
186 x_i and goal pose x_i^* for each object, we leverage a planning-based rearrangement policy to solve the
187 task. Our rearrangement policy navigates to each object in succession and transports them to their
188 goal location. By accurately mapping with a search-based policy, inferring the rearrangement goal,
189 and planning towards the goal, our method in Algorithm 1 efficiently solves visual rearrangement.

Table 1: Evaluation on the 2022 AI2-THOR 2-Phase Rearrangement Challenge. Our method attains state-of-the-art performance on this challenge, outperforming prior work by 875% *%Fixed Strict*. Results are averaged over 1000 rearrangement tasks in each of the 2022 validation set and 2022 test set. Higher is better. A *Success* of 100.0 indicates all objects are successfully rearranged if none are misplaced and 0.0 otherwise. The metric *%Fixed Strict* is more lenient, equal to the percent of objects that are successfully rearranged if none are newly misplaced, and 0.0 otherwise.

| Method | Validation | | Test | |
|--------------------------------|---------------|-------------|--------------------|--------------------|
| | %Fixed Strict | Success | %Fixed Strict | Success |
| VRR + Map [Weihs et al., 2021] | 1.18 | 0.40 | 0.53 | 0.00 |
| CSR [Gadre et al., 2022] | 3.30 | 1.20 | 1.90 | 0.40 |
| Ours w/o Semantic Search | 15.77 | 4.30 | +795% 15.11 | +900% 3.60 |
| Ours | 17.47 | 6.30 | +875% 16.62 | +1158% 4.63 |

190 4 Experiments

191 We presented a modular approach for rearrangement. In this section we evaluate our approach and
 192 show its effectiveness. We first evaluate our approach on the AI2-THOR Rearrangement Challenge
 193 Weihs et al. [2021] and show our approach leads to an improvement of 14.72 absolute percentage
 194 points over current work, detailed in Subsection 4.1. This benchmark tests an agent’s ability to
 195 rearrange rooms to a desired object goal configuration, and is a suitable choice for measuring visual
 196 rearrangement performance. Next, we show the importance of each proposed component, and
 197 demonstrate in Subsection 4.2 our voxel-based map and search-based policy exhibit large potential
 198 gains as more performant models for perception and search are developed in the future. Finally,
 199 we show in Subsection 4.3 our approach is robust to the quality of object detections and budget for
 200 exploration. Our experiments show our method is robust and effective at visual rearrangement.

201 **Description of the benchmark.** In this benchmark, the goal is to rearrange up to five objects to
 202 a desired state, defined in terms of object locations and openness. The challenge is based on the
 203 RoomR [Weihs et al., 2021] dataset that consists of a training set with 80 rooms and 4000 tasks,
 204 validation set with 20 rooms and 1000 tasks, and a test set with 20 rooms and 1000 tasks. We consider
 205 a two-phase setting where an agent observes the goal configuration of the scene during an initial
 206 *Walkthrough Phase*. The scene is then shuffled, and the agent is tasked with rearranging objects back
 207 to their goal configuration during a second *Unshuffle Phase*. This two-phase rearrangement setting is
 208 challenging because it requires the agent to remember the scene layout from the *Walkthrough Phase*,
 209 to identify the rearrangement goal. Goals are internally represented by a set of valid object poses
 210 $S^* \subset (\mathcal{R}^3 \times SO(3)) \times (\mathcal{R}^3 \times SO(3)) \dots \times (\mathcal{R}^3 \times SO(3))$, but the agent does not observe S^* directly.
 211 At every time step t during either phase, the agent observes a geocentric pose x_t , an egocentric RGB
 212 image I_t , and an egocentric depth image D_t . The rearrangement goal is specified indirectly via
 213 observations of the scene layout during the *Walkthrough Phase*. During training, additional metadata
 214 is available such as ground-truth semantic labels, but during evaluation only the allowed observations
 215 x_t , I_t and D_t can be used. Once both the *Walkthrough Phase* and *Unshuffle Phase* are complete, we
 216 measure performance using the *%Fixed Strict* and *Success* metrics described in Section 3.

217 4.1 Effectiveness At Visual Rearrangement

218 The goal of this subsection is to evaluate the effectiveness of our method at visual rearrangement.
 219 We leverage the RoomR [Weihs et al., 2021] dataset and evaluate our method on the two-phase rear-
 220 rangement challenge. We report performance in Table 1 and show an improvement in *%Fixed Strict*
 221 from 1.9 to 15.11 over the current state-of-the-art method, namely Continuous Scene Representations
 222 (CSR) [Gadre et al., 2022]. These results show our method is more effective than prior work at visual
 223 rearrangement, leading to a relative improvement of 875% over current work. Our success of 4.63%
 224 on the test set indicates our method solves 46 / 1000 tasks, whereas the best existing approach, CSR,
 225 solves 4 / 1000 tasks. Furthermore, our method correctly rearranges 499 / 3004 objects in the test set,
 226 while the best existing approach, CSR, rearranges only 57 / 3004 objects in the test set.

227 The results in Table 1 support two conclusions. First, 3D Mapping is a helpful inductive bias. Ours is
 228 currently the only method on the challenge to leverage 3D Mapping for identifying rearrangement

Table 2: Ablation of the importance of each component of our method. Our method produces significant gains as perception and search models become more accurate. Results are averaged over 1000 rearrangement tasks in each of the 2022 validation set and 2022 test set. As in Table 1, higher is better, and a *Success* of 100.0 indicates all objects are successfully rearranged if none are misplaced and 0.0 otherwise. Our results show that as perception and search models continue to improve with future research, we have an *out-of-the-box* improvement of 34.73 *Success* on the test set.

| Method | Validation | | Test | |
|---------------------------|---------------|--------------|---------------------|---------------------|
| | %Fixed Strict | Success | %Fixed Strict | Success |
| CSR + GT T | 3.80 | 1.30 | 2.10 | 0.70 |
| CSR + GT BT | 7.90 | 3.00 | 5.90 | 2.20 |
| CSR + GT MBT | 26.00 | 8.80 | 27.00 | 10.00 |
| Ours + GT Semantic Search | 21.24 | 7.60 | +942% 19.79 | +871% 6.10 |
| Ours + GT Segmentation | 66.66 | 45.60 | +1004% 59.29 | +1707% 37.55 |
| Ours + GT Both | 68.46 | 48.60 | +1008% 59.50 | +1742% 38.33 |

229 goals. The next best approach, CSR Gadre et al. [2022], represents the scene with a graph, where
 230 nodes encode objects, and edges encode spatial relationships. Determining which objects need
 231 to be rearranged benefits from knowing their fine-grain 3D position, which our method directly
 232 represents in our semantic maps. These results suggest an important hypothesis that our method more
 233 successfully rearranges small objects. This is an important contribution (see additional results in
 234 Subsection 4.5) because many common objects humans use are small—cutlery, plates, cups, writing
 235 implements, etc. Agents helpful to humans must successfully handle these small objects.

236 4.2 Component Ablation

237 The goal of this experiment is to determine the importance of each component to our method. We
 238 consider a series of ablations in Table 2 that replace different components of our method with ground
 239 truth predictions. We first consider *Ours + GT Semantic Search*, where we substitute the predictions
 240 of our search-based policy π_θ with the ground truth locations of objects that need to be rearranged.
 241 We also consider *Ours + GT Segmentation*, where we substitute the predictions of Mask R-CNN [He
 242 et al., 2017] with ground truth semantic segmentation labels. The final ablation in the table *Ours +*
 243 *GT Both* includes both substitutions at once. In addition to reporting our performance, we reference
 244 the performance of CSR [Gadre et al., 2022] in a similar set of ablations. We consider *CSR + GT T*
 245 which uses expert trajectories that observe all objects needing to be rearranged, *CSR + GT BT* which
 246 also uses ground truth object detection labels, and *CSR + GT MBT* which additionally uses ground
 247 truth object instance pairs between the *Walkthrough Phase* and the *Unshuffle Phase*. Table 2 shows
 248 our method produces a better *out-of-the-box* improvement in all metrics as the perception and search
 249 components become more accurate, suggesting both components are important.

250 Table 2 demonstrates our method produces significant gains when paired with accurate semantic
 251 search and accurate semantic segmentation. When using ground-truth semantic segmentation labels
 252 and ground-truth search locations, our method attains an improvement of 35.35 absolute percentage
 253 points in *Success* compared to existing work given access to the same experts. *CSR + GT BT* makes
 254 the same assumptions as our method with both components replaced with ground-truth, and is used
 255 to compute this improvement margin. When prior work is given the *additional* accommodation of
 256 ground-truth object instance pairs between the two environment phases, *CSR + GT MBT*, our method
 257 maintains an improvement of 27.55 absolute *Success* points without the accommodation. These
 258 results show our method has greater room for improvement than prior work, with a %Fixed Strict
 259 32.50 absolute percentage points higher than current work. Our method’s room for improvement with
 260 more accurate perception and search models is appealing because accurate 3D vision models are an
 261 active area of research, and our method directly benefits from innovations in these models.

262 4.3 Stability Versus Perception Quality

263 In the previous sections, we evaluated our method’s effectiveness at rearrangement, and room for
 264 growth as better perception and search models are developed. This direct benefit from improvements
 265 in perception quality resulting from better models is desirable, but an effective method should also

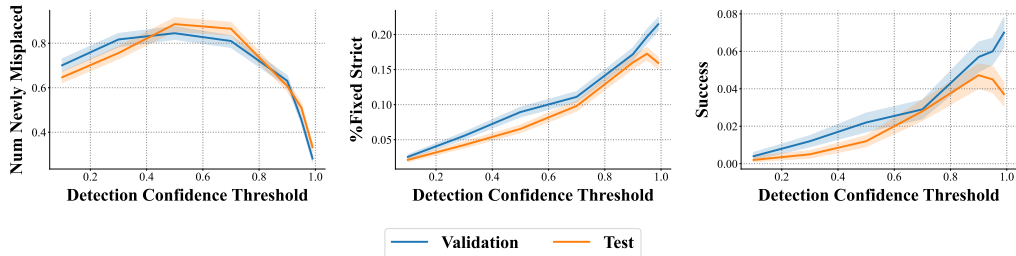


Figure 3: Rearrangement performance versus perception quality. Dark colored lines represent the average metric across 1000 tasks, and shaded regions correspond to a 68% confidence interval. Lower *Num Newly Misplaced* (left plot) is better, higher *%Fixed Strict* (center plot) and *Success* (right plot) are better. Our method improves smoothly as perception quality increases, simulated by varying the detection confidence threshold used to filter Mask R-CNN predictions detailed in Section 3.

266 be robust when perception quality is poor. In this section, we evaluate our method’s performance
 267 stability as a function of the quality of object detections. We simulate changes in object detection
 268 quality by varying the detection confidence threshold of Mask R-CNN [He et al., 2017] described in
 269 Section 3. A low threshold permits accepting detections where Mask R-CNN makes high-variance
 270 predictions, reducing the quality of detections overall. In the following experiment, we vary the
 271 detection confidence threshold on the validation and test sets of the rearrangement challenge.

272 Figure 3 shows our method is robust to small changes in perception quality. As the detection
 273 confidence increases, simulating an improvement in object detection fidelity, performance of our
 274 method smoothly increases. Peak performance with our method on the validation set is attained with
 275 a detection confidence threshold close to 0.9, which is the value we employ throughout the paper.
 276 Error bars in this experiment are computed using a 68% confidence interval with 1000 sample points,
 277 corresponding to 1000 tasks in each of the validation and test sets. The small width of error bars
 278 indicates the observed relationship between perception quality and performance most likely holds
 279 for tasks individually (not just on average), supporting the conclusion our method is robust to small
 280 changes in perception quality. We make a final observation that as perception quality increases, fewer
 281 objects are misplaced as our method more accurately infers the rearrangement goal. These results
 282 suggest our method produces consistent gains in rearrangement as perception models improve.

283 4.4 Stability Versus Exploration Budget

284 We conduct an ablation in this section to evaluate how the exploration budget affects our method.
 285 This is an important experiment because the conditions an agent faces in the real world vary, and
 286 an effective agent is robust when the budget for exploring the scene is small. We simulate a limited
 287 exploration budget by varying the amount of navigation goals used by the agent when building the
 288 semantic maps. A lower budget results in fewer time steps spent building the semantic maps, and
 289 fewer updates to voxels described in Section 3. With fewer updates, sampling goals intelligently is
 290 crucial to ensure the agent has the information necessary to infer the task rearrangement goal.

291 Figure 4 shows our method is robust when the exploration budget is small. Performance is stable
 292 when less than 5 navigation goals are proposed by our semantic search module, where no penalty
 293 in *%Fixed Strict* and *Success* can be observed. This result confirms the effectiveness of semantic
 294 search: sampled goals correspond to the locations of objects likely to need rearrangement, so even
 295 when the budget is small, these objects are already observed. The experiment also shows that as
 296 the budget decreases, fewer objects are misplaced. This is intuitive because when the budget is
 297 small, fewer objects in the environment are observed and added to the map, reducing the chance
 298 of incorrect map disagreements being proposed. Additionally, when the budget is large, the agent
 299 spends the majority of the episode in navigation, and may not have enough time left to correct map
 300 disagreements, resulting in slightly lower overall performance. These results suggest our method is
 301 effective for a variety of exploration budgets, and is robust when the budget is small.

302 4.5 Failure Modes

303 Our previous experiments showed instances where our method is effective, but an understanding
 304 of its limitations is equally important. The goal of this subsection is to identify how and why our

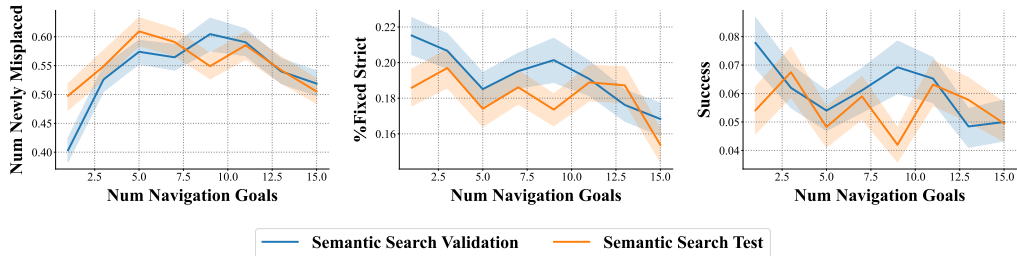


Figure 4: Rearrangement performance versus perception quality. Dark colored lines represent the average metric across 1000 tasks, and shaded regions correspond to a 68% confidence interval. Lower *Num Newly Misplaced* (left plot) is better, higher *%Fixed Strict* (center plot) and *Success* (right plot) is better. Our method improves smoothly as perception quality increases, simulated by varying the detection confidence threshold used to filter Mask R-CNN predictions detailed in Section 3.

305 method can fail. To accomplish this, we conduct an ablation to study how three indicators—object
 306 size, distance to the goal position, and amount of nearby clutter—affect our method. These capture
 307 different aspects of what makes rearrangement hard. For example, small objects can be ignored,
 308 objects distant to their goal can be easier to misplace, and objects too close to one another can be
 309 mis-detected. We measure the performance of our method with respect to these indicators in Figure 6
 310 in Appendix B, and analyze the experimental conditions when our method is less effective.

311 Our results illuminate what kinds of tasks are difficult for our method. We find experimentally that
 312 objects further from the rearrangement goal are harder for our method to successfully rearrange.
 313 Objects within 0.326 meters of the goal are correctly rearranged >30% of the time, whereas objects
 314 further than 4.157 meters from the goal are only correctly rearranged <20% of the time. One
 315 explanation for this disparity in performance could be matching object instances between phases
 316 is more difficult when those instances are further apart. Better perception models can mitigate this
 317 explanation by providing more information about object appearance that may be used to accurately
 318 pair instances. While this first observation is intuitive, our second is more surprising. We find that
 319 our method rearranges small objects as effectively as large objects, suggesting our method is robust
 320 to the size of objects it rearranges. This quality is desirable because realistic environments contain
 321 objects in a variety of sizes. Effective agents should generalize to a variety of object sizes.

322 5 Conclusion

323 We presented a simple modular approach for rearranging objects to desired visual goals. Our approach
 324 leverages a voxel-based semantic map containing objects detected by a perception model, and a
 325 semantic-search policy for efficiently locating the objects to rearrange. Our approach generalizes
 326 effectively to rearrangement goals of varying difficulties, including objects that are small in size,
 327 far from the goal, and in cluttered spaces. Furthermore, our approach is efficient, performing well
 328 even with a small exploration budget. Our experimental evaluation shows our approach improves
 329 over current work in rearrangement by 14.7 absolute percentage points, and continues to improve
 330 smoothly as better models are developed and the quality of object detections increases. Our results
 331 confirm the efficacy of active perceptual mapping for the rearrangement setting, and motivate several
 332 future directions that can expand the flexibility and generalization of the method.

333 One promising future direction is improving the map representation of objects. One limitation of
 334 the rearrangement setting in this work is that objects only have simple states: position, orientation,
 335 and openness. Real objects are complex and have states that may change over time, potentially from
 336 interactions not involving the agent. Investigating tasks that require modelling these dynamic objects
 337 in the map is an emerging topic that can benefit from new benchmarks and methods. A second
 338 promising future direction is using an agent’s experience to improve its perception. Feedback from
 339 the environment, including instructions, rewards, and transition dynamics, provides rich information
 340 about how to improve perception when true labels may be difficult to acquire. Investigating how to
 341 leverage all sources of feedback available to an agent is a useful research topic that may unlock better
 342 generalization for embodied agents in dynamic environments.

References

- 343
- 344 P. Anderson, A. Chang, D. S. Chaplot, A. Dosovitskiy, S. Gupta, V. Koltun, J. Kosecka, J. Ma-
345 lik, R. Mottaghi, M. Savva, and A. R. Zamir. On evaluation of embodied navigation agents.
346 *arXiv:1807.06757*, 2018a.
- 347 P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson, N. Sünderhauf, I. D. Reid, S. Gould, and
348 A. van den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation
349 instructions in real environments. In *CVPR*, 2018b.
- 350 D. Batra, A. X. Chang, S. Chernova, A. J. Davison, J. Deng, V. Koltun, S. Levine, J. Malik,
351 I. Mordatch, R. Mottaghi, et al. Rearrangement: A challenge for embodied ai. *arXiv:2011.01975*,
352 2020a.
- 353 D. Batra, A. Gokaslan, A. Kembhavi, O. Maksymets, R. Mottaghi, M. Savva, A. Toshev, and
354 E. Wijmans. Objectnav revisited: On evaluation of embodied agents navigating to objects.
355 *arXiv:2006.13171*, 2020b.
- 356 O. Ben-Shahar and E. Rivlin. Practical pushing planning for rearrangement tasks. *ICRA*, 1996.
- 357 V. Blukis, C. Paxton, D. Fox, A. Garg, and Y. Artzi. A persistent spatial semantic representation for
358 high-level natural language instruction execution. In *CoRL*, 2021.
- 359 G. Bradski. The OpenCV Library. *Dr. Dobb's Journal of Software Tools*, 2000.
- 360 D. S. Chaplot, D. Gandhi, A. Gupta, and R. Salakhutdinov. Object goal navigation using goal-oriented
361 semantic exploration. In *NeurIPS*, 2020a.
- 362 D. S. Chaplot, R. Salakhutdinov, A. Gupta, and S. Gupta. Neural topological slam for visual
363 navigation. In *CVPR*, 2020b.
- 364 D. S. Chaplot, M. Dalal, S. Gupta, J. Malik, and R. Salakhutdinov. SEAL: self-supervised embodied
365 active learning using exploration and 3d consistency. In M. Ranzato, A. Beygelzimer, Y. N.
366 Dauphin, P. Liang, and J. W. Vaughan, editors, *NeurIPS*, 2021.
- 367 K. Chen, J. P. de Vicente, G. Sepulveda, F. Xia, A. Soto, M. Vazquez, and S. Savarese. A behavioral
368 approach to visual navigation with graph localization networks. In *RSS*, 2019.
- 369 N. Correll, K. Bekris, D. Berenson, O. Brock, A. Causo, K. Hauser, K. Okada, A. Rodriguez,
370 J. Romano, and P. Wurman. Analysis and observations from the first amazon picking challenge.
371 *IEEE Transactions on Automation Science and Engineering*, 2018.
- 372 A. Cosgun, T. Hermans, V. Emeli, and M. Stilman. Push planning for object placement on cluttered
373 table surfaces. In *IROS*, 2011.
- 374 A. Das, S. Datta, G. Gkioxari, S. Lee, D. Parikh, and D. Batra. Embodied Question Answering. In
375 *CVPR*, 2018.
- 376 E. W. Dijkstra. A note on two problems in connexion with graphs. *Numerische mathematik*, 1959.
- 377 J. Engel, T. Schöps, and D. Cremers. Lsd-slam: Large-scale direct monocular slam. In *ECCV*, 2014.
- 378 J. Francis, N. Kitamura, F. Labelle, X. Lu, I. Navarro, and J. Oh. Core challenges in embodied
379 vision-language planning. *arXiv:2106.13948*, 2021.
- 380 S. Y. Gadre, K. Ehsani, S. Song, and R. Mottaghi. Continuous scene representations for embodied
381 AI. *arXiv:2203.17251*, 2022.
- 382 D. Gordon, A. Kembhavi, M. Rastegari, J. Redmon, D. Fox, and A. Farhadi. IQA: Visual question
383 answering in interactive environments. In *CVPR*, 2018.
- 384 S. Gupta, J. Davidson, S. Levine, R. Sukthankar, and J. Malik. Cognitive mapping and planning for
385 visual navigation. In *CVPR*, 2017.

- 386 K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016 IEEE*
387 *Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June*
388 *27-30, 2016*, pages 770–778. IEEE Computer Society, 2016. doi: 10.1109/CVPR.2016.90. URL
389 <https://doi.org/10.1109/CVPR.2016.90>.
- 390 K. He, G. Gkioxari, P. Dollár, and R. B. Girshick. Mask R-CNN. In *ICCV*, 2017.
- 391 J. E. King, M. Cognetti, and S. S. Srinivasa. Rearrangement planning using object-centric and
392 robot-centric action spaces. In *ICRA*, 2016.
- 393 A. Krontiris and K. E. Bekris. Efficiently solving general rearrangement tasks: A fast extension
394 primitive for an incremental sampling-based planner. In *ICRA*, 2016.
- 395 H. W. Kuhn and B. Yaw. The hungarian method for the assignment problem. *Naval Res. Logist.*
396 *Quart*, 1955.
- 397 B. Kuipers and Y.-T. Byun. A robot exploration and mapping strategy based on a semantic hierarchy
398 of spatial representations. *Robotics and autonomous systems*, 1991.
- 399 Y. Labbé, S. Zagoruyko, I. Kalevatykh, I. Laptev, J. Carpentier, M. Aubry, and J. Sivic. Monte-carlo
400 tree search for efficient visually guided rearrangement planning. *IEEE Robotics and Automation*
401 *Letters*, 2020.
- 402 T. Lin, P. Dollár, R. B. Girshick, K. He, B. Hariharan, and S. J. Belongie. Feature pyramid networks
403 for object detection. In *2017 IEEE Conference on Computer Vision and Pattern Recognition,*
404 *CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 936–944. IEEE Computer Society, 2017.
405 doi: 10.1109/CVPR.2017.106. URL <https://doi.org/10.1109/CVPR.2017.106>.
- 406 S. Y. Min, D. S. Chaplot, P. Ravikumar, Y. Bisk, and R. Salakhutdinov. FILM: following instructions
407 in language with modular methods. *arXiv:2110.07342*, 2021.
- 408 A. Pashevich, C. Schmid, and C. Sun. Episodic transformer for vision-and-language navigation. In
409 *ICCV*, 2021.
- 410 A. Rosinol, A. Violette, M. Abate, N. Hughes, Y. Chang, J. Shi, A. Gupta, and L. Carlone. Kimera:
411 from SLAM to spatial perception with 3D dynamic scene graphs. *Intl. J. of Robotics Research*,
412 2021.
- 413 M. Rünz and L. Agapito. Co-fusion: Real-time segmentation, tracking and fusion of multiple objects.
414 In *ICRA*, 2017.
- 415 M. Shridhar, J. Thomason, D. Gordon, Y. Bisk, W. Han, R. Mottaghi, L. Zettlemoyer, and D. Fox.
416 ALFRED: A Benchmark for Interpreting Grounded Instructions for Everyday Tasks. In *CVPR*,
417 2020.
- 418 K. P. Singh, S. Bhambri, B. Kim, R. Mottaghi, and J. Choi. Factorizing perception and policy for
419 interactive instruction following. In *ICCV*. IEEE, 2021.
- 420 M. Stilman, J.-U. Schamburek, J. Kuffner, and T. Asfour. Manipulation planning among movable
421 obstacles. In *ICRA*, 2007.
- 422 S. Thrun. Robotic mapping: A survey. *Exploring artificial intelligence in the new millennium*, 2002.
- 423 S. Thrun, D. Fox, W. Burgard, and F. Dellaert. Robust monte carlo localization for mobile robots.
424 *Artificial Intelligence*, 2001.
- 425 L. Weihs, M. Deitke, A. Kembhavi, and R. Mottaghi. Visual room rearrangement. In *CVPR*, 2021.
- 426 E. Wijmans, A. Kadian, A. S. Morcos, S. Lee, I. Essa, D. Parikh, M. Savva, and D. Batra. Dd-ppo:
427 Learning near-perfect pointgoal navigators from 2.5 billion frames. In *ICLR*, 2020.
- 428 Y.-S. Wong, C. Li, M. Niessner, and N. J. Mitra. Rigidfusion: Rgb-d scene reconstruction with
429 rigidly-moving objects. *Computer Graphics Forum*, 2021.

- 430 Y. Wu, A. Kirillov, F. Massa, W.-Y. Lo, and R. Girshick. Detectron2. [https://github.com/](https://github.com/facebookresearch/detectron2)
431 [facebookresearch/detectron2](https://github.com/facebookresearch/detectron2), 2019a.
- 432 Y. Wu, Y. Wu, A. Tamar, S. Russell, G. Gkioxari, and Y. Tian. Bayesian relational memory for
433 semantic visual navigation. In *ICCV*, 2019b.
- 434 W. Yang, X. Wang, A. Farhadi, A. Gupta, and R. Mottaghi. Visual semantic navigation using scene
435 priors. In *ICLR*, 2019.
- 436 J. Ye, D. Batra, A. Das, and E. Wijmans. Auxiliary tasks and exploration enable objectnav.
437 *arXiv:2104.04112*, 2021.
- 438 W. Yuan, J. A. Stork, D. Kragic, M. Y. Wang, and K. Hang. Rearrangement with nonprehensile
439 manipulation using deep reinforcement learning. In *ICRA*, 2018.

440 Checklist

- 441 1. For all authors...
- 442 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
443 contributions and scope? [Yes] See Sections 3 and 4.1
- 444 (b) Did you describe the limitations of your work? [Yes] See Section 4.5
- 445 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 446 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
447 them? [Yes]
- 448 2. If you are including theoretical results...
- 449 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 450 (b) Did you include complete proofs of all theoretical results? [N/A]
- 451 3. If you ran experiments...
- 452 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
453 mental results (either in the supplemental material or as a URL)? [Yes] Code will be
454 released upon acceptance. See Algorithm 1 for pseudo-code. Section 4 contains notes
455 on existing data set RoomR for AI2-THOR Rearrangement Challenge. Details about
456 hyperparameters and tuning are given in Appendix E.
- 457 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
458 were chosen)? [Yes] Details about hyperparameter tuning and training details are given
459 in Appendix E.
- 460 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
461 ments multiple times)? [Yes] While the official challenge leaderboard (performance
462 reported in the Table 1 and Table 2) does not expose error bars, we have conducted an
463 additional performance comparison in Appendix C that includes error bars.
- 464 (d) Did you include the total amount of compute and the type of resources used (e.g., type
465 of GPUs, internal cluster, or cloud provider)? [Yes] Details about necessary compute
466 are available in Appendix D.
- 467 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 468 (a) If your work uses existing assets, did you cite the creators? [Yes] See Section 4
- 469 (b) Did you mention the license of the assets? [No] The original dataset RoomR was re-
470 leased under Apache License 2.0. See <https://github.com/allenai/ai2thor-rearrangement>
- 471 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
472 No new data.
- 473 (d) Did you discuss whether and how consent was obtained from people whose data you’re
474 using/curating? [No] RoomR was released under Apache 2.0 license.
- 475 (e) Did you discuss whether the data you are using/curating contains personally identifiable
476 information or offensive content? [N/A]
- 477 5. If you used crowdsourcing or conducted research with human subjects...

- 478 (a) Did you include the full text of instructions given to participants and screenshots, if
479 applicable? [N/A]
- 480 (b) Did you describe any potential participant risks, with links to Institutional Review
481 Board (IRB) approvals, if applicable? [N/A]
- 482 (c) Did you include the estimated hourly wage paid to participants and the total amount
483 spent on participant compensation? [N/A]

(Supplementary Material)

A Simple Approach for Visual Rearrangement: 3D Mapping and Semantic Search

484 In this appendix we include the following supporting experiments and visualizations:

- 485 A. We begin this appendix by presenting the performance of our map disagreement detec-
486 tion module for each object category. We find that our method effectively detects map
487 disagreements for both small and large objects, and is therefore robust to object size.
- 488 B. We then present a performance breakdown of our method for object size, distance to goal,
489 and amount of clutter, and find that our method is less effective when objects are further
490 from the goal or when nearby objects are closer together.
- 491 C. We report confidence intervals for our method’s performance on the rearrangement challenge.
- 492 D. Finally, we outline the compute infrastructure needed to reproduce our experiments.
- 493 E. We list the hyperparameters used in our paper.
- 494 F. We categorize why our method can fail and provide a qualitative example.

495 The official code for our method will be released at publication.

496 A Object Type Versus Detection Accuracy

497 In this section, we visualize the relationship between the performance of our map disagreement
498 detection module, detailed in Section 3, and the category of objects to be rearranged. For each of
499 1000 tasks in the validation set and test set of RoomR [Weihs et al., 2021], we record which object
500 categories are detected as needing to be rearranged, and log the ground truth list of object categories
501 that need to be rearranged. For each object, we calculate precision as the proportion of objects per
502 category that were correctly identified as map disagreements out of all predicted map disagreements.
503 Similarly, we calculate recall as the proportion of correctly identified as map disagreements out of
504 all ground-truth map disagreements. Each bar in Figure 5 represents a 68% confidence interval of
505 precision and recall over 1000 tasks per dataset split. The experiment shows that our method is
506 robust to the size of objects that it rearranges because small objects such as the *SoapBar*, *CellPhone*,
507 *CreditCard*, and *DishSponge* have comparable accuracy to large objects in Figure 5.

508 B Performance Analysis

509 This section extends Section 4.5 with an experiment to show potential failure modes. We consider
510 three failure modes: (1) object size, (2) object distance to the goal, and (3) closest object in the
511 same class. These indicators are visualized in Figure 6 against *%Fixed*. Our experiment suggests
512 our method is robust to the size of objects, shown by the lack of a global trend in the left plot in
513 Figure 6, and confirmed by Appendix A. Additionally, the experiment shows that objects further from
514 the rearrangement goal are solved less frequently (middle plot), which is intuitive. Instances that
515 have been shuffled to faraway locations in the scene may require longer exploration to find, and may
516 be more difficult for our map disagreement detection module to match. A final conclusion we can
517 draw from this experiment is that our method can fail when object instances are too close together.
518 This is shown in the right plot in Figure 6 by the steep drop in performance when objects in the
519 same category are < 1 meter apart. In this situation, our semantic mapping module can incorrectly
520 detect two nearby objects as a single object, which prevents their successful rearrangement. For each
521 of these potential failure modes, better perception and mapping approaches that more accurately
522 describe object locations and appearance can improve fidelity of our method and reduce failure.

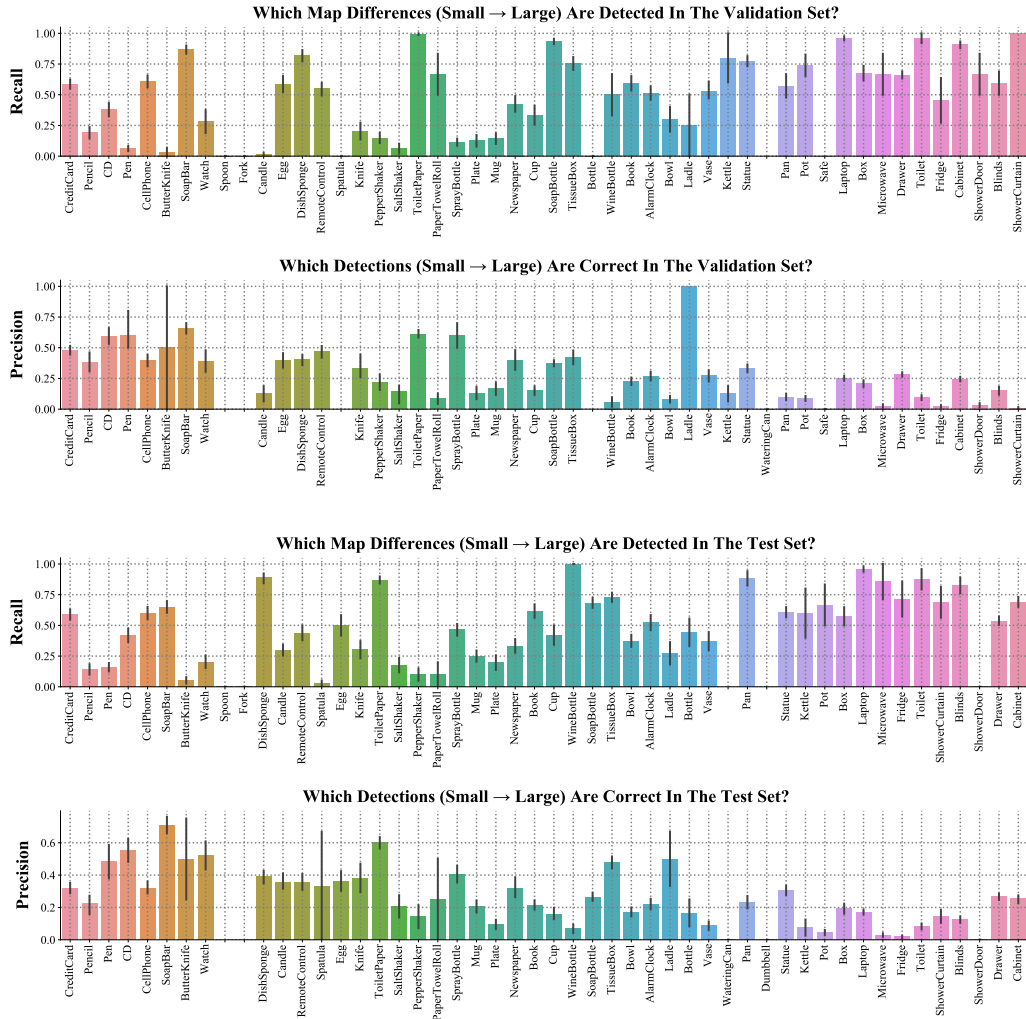


Figure 5: Performance breakdown on the validation and test sets for various types of objects. The height of bars corresponds to the sample mean of precision or recall for our map disagreement detection module. Error bars show a 68% confidence interval for each kind of object. The top two plots correspond to precision and recall on the validation set, while the bottom two plots correspond to precision and recall on the test set. Object categories are shown on the x-axis, and are ordered in ascending order of size. The experiment shows our method is robust to size, with small objects at the left end of the plots having comparable accuracy to large objects at the right end of the plots.

523 C Performance Confidence Intervals

524 We report 68% confidence intervals in Table 3 to supplement our evaluation in Section 4.1 and Sec-
 525 tion 4.2. We calculate intervals using 1000 tasks from the validation and test sets of the RoomR [Weihs
 526 et al., 2021] dataset, and report the mean followed by \pm interval width. Note that the official rearrange-
 527 ment challenge leaderboard does not expose confidence intervals, nor the sample-wise performance
 528 needed to calculate them. Due to this, we are unable to compute confidence intervals of the baselines
 529 VRR [Weihs et al., 2021] and CSR [Gadre et al., 2022] at this time. These additional results show
 530 that our improvements over prior work significantly exceed the 68% confidence interval.

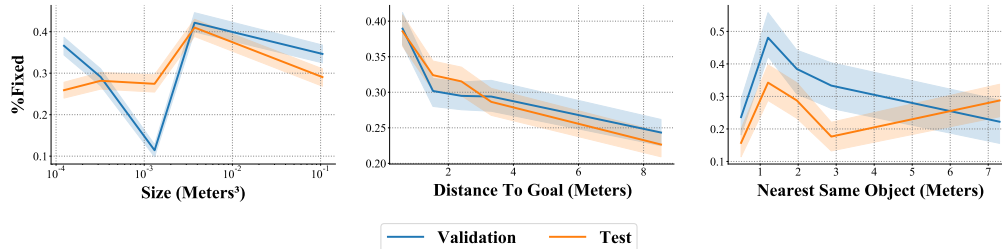


Figure 6: Performance of various ablations for different *Size (Meters³)*, *Distance To Goal (Meters)*, and *Nearest Same Object (Meters)*. These indicators measure properties of objects that make rearrangement hard. Colored lines represent the average performance over 1000 tasks in each dataset split. Error bars represent a 68% confidence interval over those same 1000 sample points. The experiment shows our method can fail when objects of the same class are too close together (right plot), and when objects are too far from the goal location, typically >4.157 meters (center plot).

Table 3: Confidence intervals for our method on the AI2-THOR rearrangement challenge. Intervals are calculated from 1000 sample points from RoomR [Weihs et al., 2021] validation and test sets. We report performance starting with the sample mean, followed by \pm a 68% confidence interval width. Our improvements over prior work significantly exceed the 68% confidence interval, which suggests that our improvements are significant and our method performs consistently well.

| Method | Validation | | Test | |
|---------------------------|------------------|------------------|------------------|------------------|
| | %Fixed Strict | Success | %Fixed Strict | Success |
| Ours w/o Semantic Search | 15.77 \pm 0.85 | 4.30 \pm 0.63 | 15.11 \pm 0.84 | 3.60 \pm 0.58 |
| Ours | 17.47 \pm 0.92 | 6.30 \pm 0.76 | 16.62 \pm 0.89 | 4.63 \pm 0.67 |
| Ours + GT Semantic Search | 21.24 \pm 0.99 | 7.60 \pm 0.83 | 19.79 \pm 0.96 | 6.10 \pm 0.75 |
| Ours + GT Segmentation | 66.66 \pm 1.21 | 45.60 \pm 1.57 | 59.29 \pm 1.26 | 37.55 \pm 1.53 |
| Ours + GT Both | 68.46 \pm 1.20 | 48.60 \pm 1.57 | 59.50 \pm 1.31 | 38.33 \pm 1.57 |

531 D Required Compute

532 The goal of this section is to outline the amount of compute required to replicate our experiments.
 533 We will describe the amount of compute required for (1) training Mask R-CNN, (2) training a
 534 semantic search policy $\pi_\theta(x|m_i)$, and (3) benchmarking the agent on the rearrangement challenge.
 535 For training Mask R-CNN, a dataset of 2 million images with instance segmentation labels were
 536 collected from the THOR simulator using the training split of the RoomR [Weihs et al., 2021] dataset.
 537 We then used Detectron2 [Wu et al., 2019a] with default hyperparameters to train Mask R-CNN with
 538 a ResNet50 [He et al., 2016] Feature Pyramid Network backbone [Lin et al., 2017]. We trained our
 539 Mask R-CNN for five epochs using a DGX with eight Nvidia 32GB v100 GPUS for 48 hours. Our
 540 semantic search policy requires significantly less compute: completing 15 epochs on a dataset of
 541 8000 semantic maps annotated with an expert search distribution in nine hours on a single Nvidia
 542 12GB 3080ti GPU. Evaluating our method on the AI2-THOR rearrangement challenge requires 40
 543 GPU-hours with a 2080ti GPU or equivalent. In practice, we parallelize evaluation across 32 GPUs,
 544 which results in an evaluation time of 1.25 hours for each of the validation and test sets.

545 E Hyperparameters

546 We provide a list of hyperparameters and their values in Table 4. These hyperparameters are held
 547 constant throughout the paper, except in ablations that study the sensitivity of our method to them,
 548 such as Section 4.3. Our ablations show our method is robust to these hyperparameters.

Table 4: Hyperparameters used by our approach for all rearrangement tasks.

| Hyperparameter | Value |
|--------------------------------------|--------------|
| voxel size | 0.05 meters |
| map height H | 384 |
| map width W | 384 |
| map depth D | 96 |
| classes C | 54 |
| detection confidence threshold | 0.9 |
| rearrangement distance threshold | 0.05 meters |
| expert search distribution σ | 0.75 meters |
| π_θ convolution hidden size | 64 |
| π_θ convolution kernel size | 3×3 |
| π_θ layers | 5 |
| π_θ activation function | ReLU |
| π_θ optimizer | Adam |
| π_θ learning rate | 0.0003 |
| π_θ batch size | 8 |
| π_θ epochs | 15 |
| π_θ dataset size | 8000 |

549 F Reasons For Task Failures

550 This section explores the reasons why certain tasks in the validation and test sets are not solved by
 551 our method. We consider four reasons for task failures that cover all possible outcomes: (1) the agent
 552 correctly predicts which objects need to be moved where, but fails to rearrange at least one object, (2)
 553 the agent incorrectly predicts an object needs to be rearranged that doesn't, (3) the agent runs out
 554 of time, and (4) the agent misses at least one object that needs to be rearranged. We visualize the
 555 proportion of failed tasks for each category in Figure 7. We find that our method with ground truth
 556 perception and search (*Ours + GT SSS*) tends to fail to rearrange objects after correctly identifying
 557 which objects need to be rearranged. In contrast, the largest reason for failure for our method (*Ours*)
 558 is the agent running out of time, followed by rearranging incorrect objects. This suggests the largest
 559 potential gains for our method arise from improving the speed and fidelity of map building, whereas,
 560 the optimality of the rearrangement policy becomes the bottleneck once a perfect map is available.

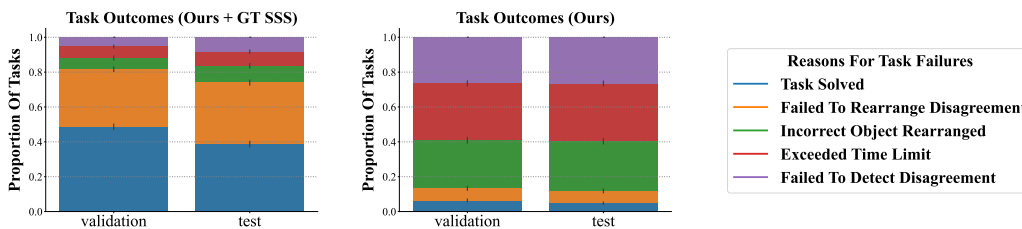


Figure 7: Categorization of the reasons why our method fails to solve tasks. The proportion of tasks that are solved (shown in blue) or fail due to one of four reasons (orange, green, red, purple) is shown for different ablations of our method. The height per bar corresponds to the proportion of tasks in the validation or test set in each category, and error bars indicate a 68% confidence interval. This experiment shows the largest reason for failure is a result of mapping errors. In the right plot, the agent fails most frequently by rearranging the wrong object, and by running out of time, which can result from imperfect semantic maps. In contrast, once perfect maps are available in the left plot, the largest source of errors are due to an imperfect planning-based rearrangement policy instead.



Figure 8: Qualitative example for why rearranging the correct object can fail. In this task, the agent correctly predicts the *ToiletPaper* needs to be rearranged, but fails to place the *ToiletPaper* in the correct location. The rightmost image shows the goal is located on the floor, but the agent mistakenly places the *ToiletPaper* on the bathtub instead, shown in the second image from the right.