# Object Manipulation via Visual Target Localization

Kiana Ehsani[1]    Ali Farhadi[2]    Aniruddha Kembhavi[1,2]    Roozbeh Mottaghi[1,2]

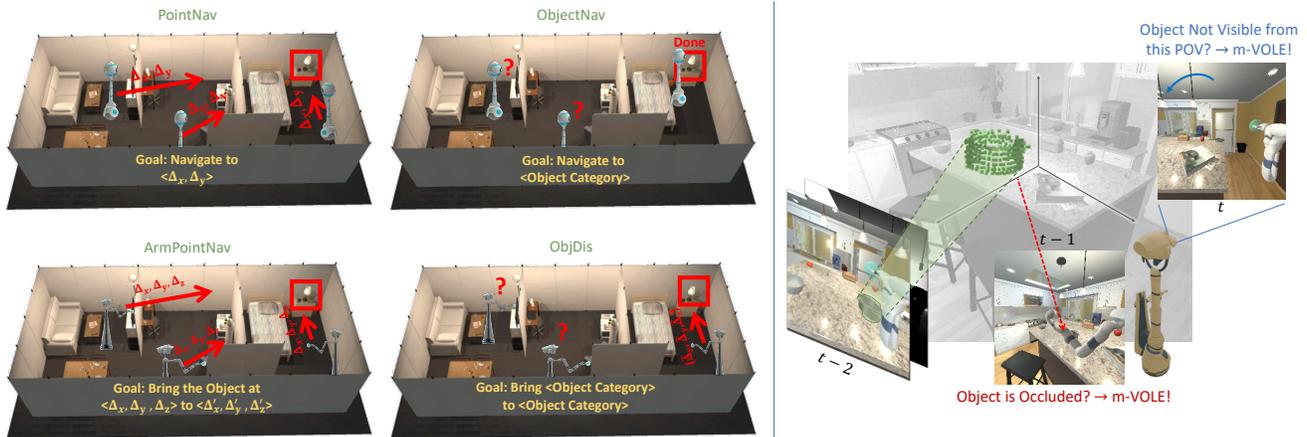[1] Allen Institute for AI   [2] University of Washington

Figure 1. We propose to solve a manipulation task, Object Displacement (ObjDis), where the goal is to bring a source object towards a destination object (e.g., bring a bowl to sink). In previous popular Embodied AI tasks shown in the left panel, either the goal is defined by the relative 3D coordinates (PointNav [38] and ArmPointNav [11]) or there is no manipulation involved (ObjectNav [4]). In contrast, we estimate the goal location from visual observations and manipulate objects across the scenes. The right panel shows an example that the agent robustly estimates the relative object location despite the occlusion by the arm and being out of view.

## Abstract

*Object manipulation is a critical skill required for Embodied AI agents interacting with the world around them. Training agents to manipulate objects, poses many challenges. These include occlusion of the target object by the agent's arm, noisy object detection and localization, and the target frequently going out of view as the agent moves around in the scene. We propose Manipulation via Visual Object Location Estimation (m-VOLE), an approach that explores the environment in search for target objects, computes their 3D coordinates once they are located, and then continues to estimate their 3D locations even when the objects are not visible, thus robustly aiding the task of manipulating these objects throughout the episode. Our evaluations show a massive 3x improvement in success rate over a model that has access to the same sensory suite but is trained without the object location estimator, and our analysis shows that our agent is robust to noise in depth perception and agent localization. Importantly, our proposed approach relaxes several assumptions about idealized localization and perception that are commonly employed by recent works in embodied AI – an important step towards training agents for object manipulation in the real world.*

## 1. Introduction

In recent years the computer vision community has made steady progress on a variety of Embodied AI tasks including navigation [1, 4, 6, 57], object manipulation [11, 43, 44, 48] and language-based tasks [2, 9, 14, 41] within interactive worlds in simulation [12, 24, 38, 49]. Performances of state-of-the-art systems on these tasks vary greatly depending on the complexity of the task, the assumptions made about the agent and environment, and the sensors employed by the agent. While robust and reliable methods have emerged for some of these tasks, developing generalizable and scalable solutions remains a topic of research. A challenging problem in this domain is visual object manipulation, a critical skill that enables the agents to interact with objects and change the world around them. Avoiding the collision of the arm with other objects in the scene, inferring the state

of the scene using noisy visual observations, and planning an efficient path for the agent and its arm towards objects are a few of many interesting challenges in this domain.

One of the main obstacles in training capable and generalizable agents for visual object manipulation is the sparsity of training signals. Early works in visual navigation in the Embodied AI (EAI) research community (e.g., [57]) suffered from this training signal problem. To alleviate this issue, the EAI community provided the agent with a powerful suite of sensors. The first incarnations of the Point Goal Navigation (PointNav) task (navigating to a specified X-Y coordinate within an environment) relied on perfect sensory information that included GPS localization for the target along with compass and GPS sensors for the agent [38]. As a result, the agent was able to access accurate relative coordinates of its goal at every time step, which acted as a dense supervisory signal during training and inference. Under these assumptions, models with minimal inductive biases achieve near-perfect accuracy [46] in unseen test environments given enough training. Mirroring this success story, researchers have also been able to train effective agents for visual object manipulation [11].

While these are promising steps towards building embodied agents, the strong sensory assumptions employed by these models limit their applicability in the real world. Their task definition requires specifying the target via 3D coordinates, which can be extremely hard to determine, if not impossible, in indoor environments. Moreover, they rely on a perfect compass and GPS sensory information, which enables the agent to localize itself with respect to the target. In this paper, we take a step closer to a more realistic task definition for object manipulation by specifying the goal via a representative image of a particular object category (instead of their 3D coordinates) and propose a method to localize the target object without relying on any perfect localization sensory information.

We introduce Object Displacement (ObjDis), the task of bringing an object to a target location (e.g., *bring a bowl to the sink*). This involves searching for an object, navigating to it, picking it up, and placing it at the desired target location. Figure 1(left) contrasts ObjDis with other popular navigation and manipulation tasks. We propose Manipulation via Visual Object Location Estimation (m-VOLE), a model for ObjDis that continually estimates the 3D relative location of the target to the agent via visual observations and learns a policy to perform the desired manipulation task. m-VOLE predicts a segmentation mask for the objects of interest and leverages the depth sensor to estimate the relative 3D coordinates for these objects. However, target segmentations are not always available to the agent due to several reasons including: objects may be out of view due to the agent being in a different location, objects may be occluded by the arm of the agent, and masks may be unavail-

able due to imperfect object segmentation models; and these noisy observations of the target can lead to the failure of the agent. To alleviate these issues, m-VOLE aggregates estimates for the objects' 3D coordinates over time, leverages previous estimates when the object mask is unavailable, and can seamlessly re-localize the object in its coordinate frame when the agent observes it again – rendering the model robust to noise in movement and perception. Figure 1(right) shows a schematic of the agent's object localization in the presence of occlusion.

We conduct our experiments using the ManipulaTHOR [11] framework which provides visually rich, physics-enabled environments with a variety of objects. We show that:

(a) Our model achieves 3x success rate compared to a model that is trained without target localization but with an identical set of sensory suites.

(b) Our model is robust against noise in perception and agent movements compared to the baselines.

(c) Our method allows for zero-shot manipulation of novel objects in unseen environments.

## 2. Related Works

**Robotic manipulation.** Manipulation and rearrangement of objects is a long-lasting problem in robotics [20, 22, 26]. Various approaches assume the full visibility of the environment and operate based on the assumption that a perception module provides the perfect observation of the environment [13, 19, 23, 25]. In contrast to these approaches, we focus on the visual perception problem and solve the task when the agent has partial and noisy observations. Several approaches have been proposed (e.g., [8, 27, 52, 54]) that address visual perception as well. However, the mentioned works focus on the tabletop rearrangement of objects. In contrast, we consider mobile manipulation of objects, which is a more general problem. Mobile manipulation has been explored in the literature as well. However, they typically use a single environment to develop and evaluate the models [5, 33, 34, 42]. One of the important problems we address in this paper is the generalization to unseen environments and configurations.

**Manipulation in virtual environments.** Recently, various works have addressed the problem of object manipulation and rearrangement in virtual environments. These works typically abstract away grasping and, unlike the works mentioned above, mostly focus on visual perception for manipulation, learning-based planning, and generalization to unseen environments and objects. Robosuite [58], RL-Bench [21] and Meta-world [53] provide a set of benchmarks for tabletop manipulation. ManiSkill benchmark [31] built upon the Sapien framework [50] is designed to benchmark manipulation skills over diverse objects using a static

robotic arm. In contrast to these works, we focus on object displacement (i.e., navigation and manipulation) in unseen scenes (e.g., a kitchen). [48] and [43] propose a set of tasks in the iGibson [39] and Habitat 2.0 [43] frameworks. However, the same environment and objects are used for train and test. In contrast, our focus is on generalization to unseen environments. ManipulaTHOR [11] is an object manipulation framework that employs a robotic arm and introduces a task and benchmark for mobile manipulation. Their task highly depends on various accurate sensory information, which renders the task unrealistic. In contrast to this work, we relax most supervisory signals to better mimic real-world scenarios.

**Relaxing supervisory sensors.** There have been some attempts to relax the assumptions about perfect sensory information using visual odometry for the navigation task [10,56]. While these are effective approaches for PointGoal navigation, the same approach does not apply to manipulation as they still require the goal's location in the agent's initial coordinate frame. In contrast, we have no access to the target or goal location at any time step.

**Embodied interactive tasks.** Navigation is a crucial skill for embodied agents. Various recent techniques have addressed the navigation problem [1,4,6,16,35,37,46,47,51, 57]. These tasks mostly assume the environments are static, and objects do not move. We consider joint manipulation and navigation, which poses different challenges. [49, 55] assume objects can move during navigation, but the set of manipulation actions is restricted to push actions. [2,30,41] propose instruction following to navigate and manipulate objects. They are designed either for static environments [2] or abstract object manipulation [30, 41] (e.g., objects snap to the agent). Works such as [28, 29] use auxiliary tasks to overcome the issues related to sparse training signals for navigation. We focus on manipulation, which deals with a different set of challenges. [3, 44] propose room rearrangement tasks using mobile robots. Similar to our paper, it involves navigation and displacement of objects. However, their object manipulation is unrealistic in that they assume a magic pointer abstraction, i.e., no arm movement is involved. In contrast, we manipulate objects using an arm, which introduces unique challenges such as planning the motion of the physical arm and avoiding collisions of the arm and object in hand with the rest of the environment.

## 3. Object Displacement

We introduce the task of Object Displacement (ObjDis), which requires an agent to navigate and manipulate objects in visually rich, cluttered environments. Given two object references, a source object $O_S$ and a destination object $O_D$, the goal is to locate and pick up $O_S$ and move it to the proximity of $O_D$ (e.g., *bring an egg to a pot*). The objects of interest are specified to the agent via query images of

an instance of each of the categories. Referring to objects via query images as opposed to object names enables us to task the agent with manipulating objects that it has not been trained for. Note that the query images *do not match* the appearance of objects within the scene but are canonical object images (we obtain the images from simulation, but they can also be obtained via other sources such as an image search engine). This enables the user to easily specify the task without the knowledge of object instances in the environment.

The ObjDis task consists of multiple implicit stages. To be successful, an agent must (1) explore its environment until it finds the source object, (2) navigate to it, (3) move its arm to the object so that its gripper may pick it up, (4) locate the destination object within the environment (5) navigate to this object and (6) place the object within its gripper in the proximity of the destination object. In visually rich and cluttered environments such as the ones present in AI2-THOR [24], this poses several challenges to the agent. *Firstly*, the agent must learn to move its body (navigate) as well as its arm (manipulate) effectively towards objects of interest to complete the desired task. *Secondly*, the agent must avoid collisions with other objects in the scene, which may occur with its body, its arm, or the source object once it has been picked up. *Thirdly*, the agent must be able to plan its actions over long horizons as it involves multiple objects, exploration, and manipulation. *Finally*, the agent must overcome noisy perception caused by frequent obstruction of its view due to its occluding arm, noisy depth sensors, and imperfect visual processing such as object detection.

We situate our experiments in AI2-THOR [24], a simulated set of environments built in Unity with a powerful physics engine, and adopt the agent provided by the ManipulaTHOR [11] framework. This agent has a rigid body with a 6DOF arm. The agent is capable of moving its body and arm in the environment while satisfying the kinematic constraints of the arm. Grasping is abstracted to a magnetic gripper (i.e., the object can be grasped if the gripper touches the object). While object grasping is a rich and challenging problem, using a magnetic gripper enables us to focus our efforts on other challenges, including exploration, navigation, arm manipulation, and long-horizon planning.

The action space for the agent consists of 11 actions. There are three agent movement actions (*MoveAhead*, *RotateRight*, and *RotateLeft*), two arm base movements (*MoveUp*, *MoveDown*), which move the base of the arm up and down with respect to the body of the agent, and six arm movements that move the gripper in x, y, z directions in agent's coordinate frame. Similar to [11], the agent movement actions move/rotate the agent's body by 0.2m/45 degrees, and the arm movements move the base of the arm or the gripper by 5cm. The framework uses inverse kinematic to calculate the final position of each joint.

# 4. Manipulation via Visual Object Location Estimation

One of the main challenges of the ObjDis task is finding and localizing the object of interest in the scenarios that the object is out of the field of view, is not detected, or is occluded by the arm. We propose a model to find the target, estimate its location, keep track of the location over time, and plan a path to accomplish the task.

Agents tackling the task of ObjDis accept as input ego-centric RGB and Depth observations along with query images for the source and destination objects. They must choose one of 11 actions at each step and follow a long sequence of steps to succeed. In the absence of perfect perception and localization, this task proves quite hard to learn, and models fail to generalize to new environments (as shown in Sec. 5).

## 4.1. Estimating Relative 3D Object Location

We first explain how the relative location of the object of interest is estimated. This is not possible for an unobserved object because its location remains unknown. However, once the object comes into view (while exploring the environment using the policy described later in Sec. 4.3), the agent can start this estimation. Note that this is in contrast to the ARMPOINTNAV [11] task where the groundtruth location of the target is provided at all time steps.

At time step $t$, if the desired objects are in view, one can generate their segmentation masks (we explain in Sec. 4.2 how we obtain these masks). We refer to these segmentations as $M_S$ and $M_D$ for the source and destination objects, respectively. Given a segmentation mask and the depth observation, the object's center in 3D is estimated by backprojecting the depth map within the segmentation mask into the 3D local coordinate of the agent. Formally, $\hat{d}_O = \pi(D[M_O], K)$, where $\hat{d}_O$ is the estimated center of the object $O$ in the agent's coordinate frame, $M_O$ is the segmentation mask of the object in the current observation, $D[M_O]$ is the observed depth frame masked by the object's segmentation (i.e., depth values for the visible regions of the object), $K$ is the intrinsic matrix of the camera, and $\pi$ is a function that projects pixels into the 3D space based on the depth values and camera intrinsics.

There are various sources of noise in object instance segmentation that make it challenging to have a reliable estimate of the 3D object location. First, segmentation models are far from perfect. They often have false positives and false negatives. Policies that rely on these masks or a quantity derived from these masks would not be able to produce a reliable sequence of actions. Second, the arm may obstruct the view of the agent as it moves in the scene – this causes the segmentation mask to periodically disappear. Third, the object might go outside the agent's camera frame after it is
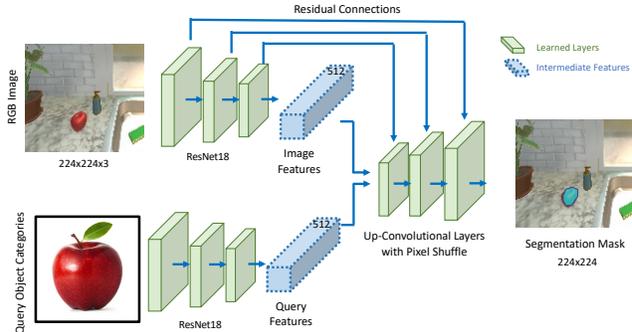


Figure 2. **Conditional Segmentation Architecture**. We use a conditional segmentation model to estimate the segmentation mask. The network receives an RGB image and a query image (representing an object category) as input and outputs the segmentation mask for an instance of that category.

observed once and as a result not visible. To overcome these issues, we propose to aggregate this information over time.

To sequentially aggregate the relative coordinates of the object, a weighted average of the previous and the current estimated 3D distance is used. At each time step (after observing the target once), m-VOLE calculates the distance $\hat{d}_O$ in the agent's coordinate frame at time $t$. However, as the agent takes action, its coordinate frame keeps shifting. At each time step that the target object is visible, the agent re-localizes the goal in its current observation frame and can readjust the coordinates. To accurately convert all the past estimates to the current agent frame, we must keep track of the agent's relative location $L_t^r$ with respect to its starting location. Since we aim to reduce the reliance on perfect sensing, we evaluate the effect of noise in the location estimation (see Sec. 5).

## 4.2. Conditional Segmentation

Estimating the 3D location of an observed object requires the agent to estimate its segmentation mask. We refer to this segmentation as 'conditional' since the goal is to obtain a mask for an object instance from the object category shown in the conditioning query image. Formulating the problem as a conditional segmentation problem as opposed to the traditional instance segmentation enables the network to focus on generating the mask for the target object[1], and the task reduces to a simpler task of category matching.

We train an auxiliary segmentation network in an offline setting, independent of the policy network. This allows us to employ a fully supervised setting resulting in efficient policy learning. This network accepts as input two images, the ego-centric RGB observation and a query object, and outputs a mask for all objects from the desired category. This reduces to a binary classification problem at each pixel in the image. The two input images are encoded via

---

[1]We use *target object* to refer either the source or destination object.

| Standard RL Reward | | Motivating Exploration | | Efficient Manipulation | |
|---|---|---|---|---|---|
| Failed Action | -0.03 | Object Observed | +1 | Pick Up $O_S$ | +5 |
| Step | -0.01 | Visit New State | +0.1 | $\Delta$ distance to obj | -$\delta$ |
| Episode Success | +10 | | | | |

Table 1. **Rewards.** In addition to standard RL reward for embodied AI, we include reward components encouraging exploration and efficient manipulation.

ResNet18 [18] models. The features are then concatenated and passed through five upconvolutional layers. The upconvolutional layers are implemented following PixelShuffle [40]. The final output is a segmentation mask of the same size as the original image (details in Figure 2). We use the Cross-Entropy loss to train this model.

We train and evaluate this model using an offline dataset collected from the train scenes in AI2-THOR [24], where we randomize object locations, textures, and scene lighting. We initialize the ResNet18 models with ImageNet pretraining and finetune them on our offline dataset. We compare this method with using a state-of-the-art instance segmentation model trained on the same data in Section 5.3.

### 4.3. Policy Network

Figure 3 shows a schematic of our proposed model m-VOLE. m-VOLE receives the RGBD observation $I_t$ at the current timestep $t$, and the query images $O_S$ and $O_D$. The conditional segmentation module estimates the segmentation masks for the target objects $M_S$ and $M_D$ and the visual object location estimator calculates the relative object coordinates $\hat{d}_o$. The policy network then uses $I_t, O_S, O_D, M_S, M_D, \hat{d}_o$ to sample actions.

The RGBD observation and segmentation masks of the target objects in the current observation $M_S, M_D$ are concatenated depth-wise and embedded using three convolutional layers. This visual encoding is combined with ResNet18 [18] (pre-trained on ImageNet [36]) features of the query images of target classes $O_S, O_D$ and the embedding of object's location $\hat{d}_o$. The resulting feature vector is provided to an LSTM, and a final linear layer generates a distribution over the possible actions. We use DD-PPO [46] to train the model. Since our method does not have access to the object's location, and the object is not necessarily initially in sight (only in 13.9% of episodes the object is initially visible to the agent), our agent is required to explore the environment until the object is found. We use rewards shown in Table 1 to encourage exploration and more efficient manipulation. In Appendix B, we ablate the effects of these rewards on the agent's performance.

## 5. Experiments

We present our experiments comparing m-VOLE with baselines (Sec. 5.1), a robustness analysis of our model with regards to noise in the agent's motions and sensors (Sec. 5.2) and ablations of m-VOLE's design choices (Sec. 5.3). Finally we evaluate m-VOLE for displacing novel objects, not used during training (Sec. 5.4).

**Baselines:** We consider the following baselines:

- **No Mask Baseline** – This model uses the RGBD observation and the query image as input and directly predicts a policy for completing the task. It does not have a segmentation module and does not have access to the agent's location relative to the starting point.

- **Mask Driven Model (MDM)** – This network uses a very similar architecture as our model with RGBD observation, query images, and the segmentation mask of the target objects as inputs. However, it does not have access to the agent's location relative to the starting point.

- **Location-aware Mask Driven Model (Loc-MDM)** – This baseline shares the same architecture as MDM. However, it also uses the agent's relative location from the start of the episode. The inputs available to Loc-MDM are the same as our proposed model (m-VOLE). The difference is that Loc-MDM uses the agent's location naively, whereas m-VOLE uses it for target localization.

In addition to these baselines, we also train the model from [11] on our task ObjDis. Note that this model uses the perfect location of the agent and the target objects at each time step, so it is not a fair comparison to our model. However, it is a useful point of reference.

- **ArmPointNav** [11] – This method is the same architecture as the one introduced in [11] with the addition of the query images as input. The network is provided with the perfect location information for the agent and source and destination objects.

**Training Details.** We train our models for 20M frames (unless otherwise specified). All the models share the same visual encoder, three convolutional layers to embed the RGBD observation. The maximum episode length is 200 frames, and the episode fails if the agent does not finish the task before the end of the episode. Using AI2-THOR [24] multi-node training, we render 500 frames per second on 4 machines, each with 4 Tesla T4 GPUs. We use Adam optimizer with a learning rate of 3e-4 and take gradient steps every 128 frames. We consider two objects being in proximity if the centers of objects are less than 20cm apart. We train our conditional segmentation model offline. We use AllenAct [45] framework.

**Dataset.** We use the APND dataset proposed in [11] to define the tasks' configurations. APND consists of 12 object categories (Apple, Bread, Tomato, Lettuce, Pot, Mug, Potato, Pan, Egg, Spatula, Cup, and SoapBottle), 30 kitchen scenes, 130 agent initial locations per scene, and 400 initial locations per object-pair per room. For each task, we
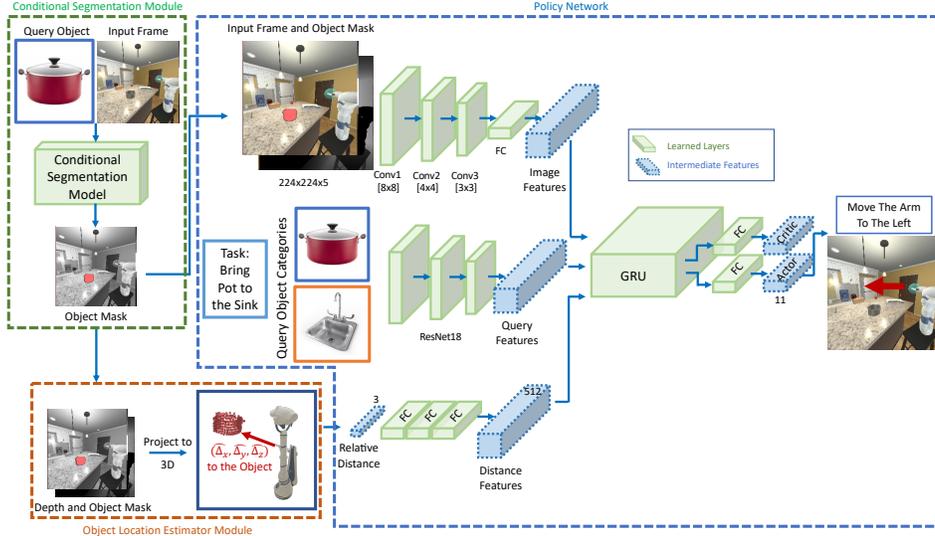
Figure 3. **Model architecture**. Our model, m-VOLE, uses the RGBD observation and a canonical query image of the object to 1) predict the *target object*'s mask (Conditional Segmentation Module), 2) estimate the relative distance of the *target object* in agent's coordinate frame (Object Location Estimator Module), and 3) predict the next action (Policy Network). Note that we use *target object* to refer to either the source or destination object.

| Model | Segmentation | Additional Input | PU | SR | SRwD |
|---|---|---|---|---|---|
| (1) No Mask Baseline | Prediction | N/A | 0.0 | 0.0 | 0.0 |
| (2) Mask Driven (MDM) | Prediction | N/A | 16.6 | 1.62 | 0.09 |
| (3) Loc-MDM | Prediction | Agent's Relative Loc | 20.3 | 3.24 | 1.08 |
| (4) m-VOLE (Ours) | Prediction | Agent's Relative Loc | **38.7** | **11.6** | **4.59** |
| (5) Mask Driven (MDM) | GT mask | N/A | 50.5 | 17.1 | 8.74 |
| (6) Loc-MDM | GT mask | Agent's Relative Loc | 57.3 | 21.6 | 9.01 |
| (7) ArmPointNav [11] | N/A | Compass + GPS | 76.6 | 58.3 | 28.5 |
| (8) m-VOLE (Ours) | GT mask | Agent's Relative Loc | **81.2** | **59.6** | **31.0** |

Table 2. **Quantitative Results.** Rows (1)-(4) present the results with the predicted segmentation masks. Rows (5)-(8) present the results for models when provided with ground truth segmentation. Our model m-VOLE outperforms all the baselines across all metrics. The No Mask Baseline is simply unable to train well, in spite of repeated attempts by us to vary hyperparameters and other design choices.

choose two objects, place them in two locations (randomly chosen from APND) and choose one as the source object $O_S$ and the other as the destination $O_D$. We also randomize the initial location of the agent. We split the 30 scenes into 20 for training, 5 for validation, and 5 for testing. We generate 132 object pairs per training scene and 1600 pairs of initial locations of objects per training scene. We select a fixed set of 1100 tasks evenly distributed across scenes and object categories for validation and test.

We also generate a dataset of images from the training scenes to train our segmentation model. Our dataset consists of more than 23K images, including 116K object masks from 120 scenes (where 80 are used for training, 20 for validation, and 20 for test sets). We ensure that the validation and test scenes used to train m-VOLE are not used to train the segmentation model. As explained in Section 3, the model requires query images as input. We collect these query images by taking a crop containing the object of in-

terest from the images of this dataset. We only use the query images collected from the training scenes of AI2-THOR, so the instances depicted in the query images do not overlap with the instances of target objects during inference.

**Metrics.** We use the same evaluation metrics as [11]:

• Success rate **(SR)** – Fraction of episodes in which the agent successfully picks up the source object ($O_S$) and brings it to the destination object ($O_D$).
• Success rate without disturbance **(SRwD)** – Fraction of episodes in which the task is successful, and the agent does not move other objects in the scene.
• Pick up success rate **(PU)** – Fraction of episodes where the agent successfully picks up the object.

### 5.1. How well does m-VOLE work?

Table 2 presents our primary quantitative findings on the test environments. Rows 1-4 provide results for m-VOLE
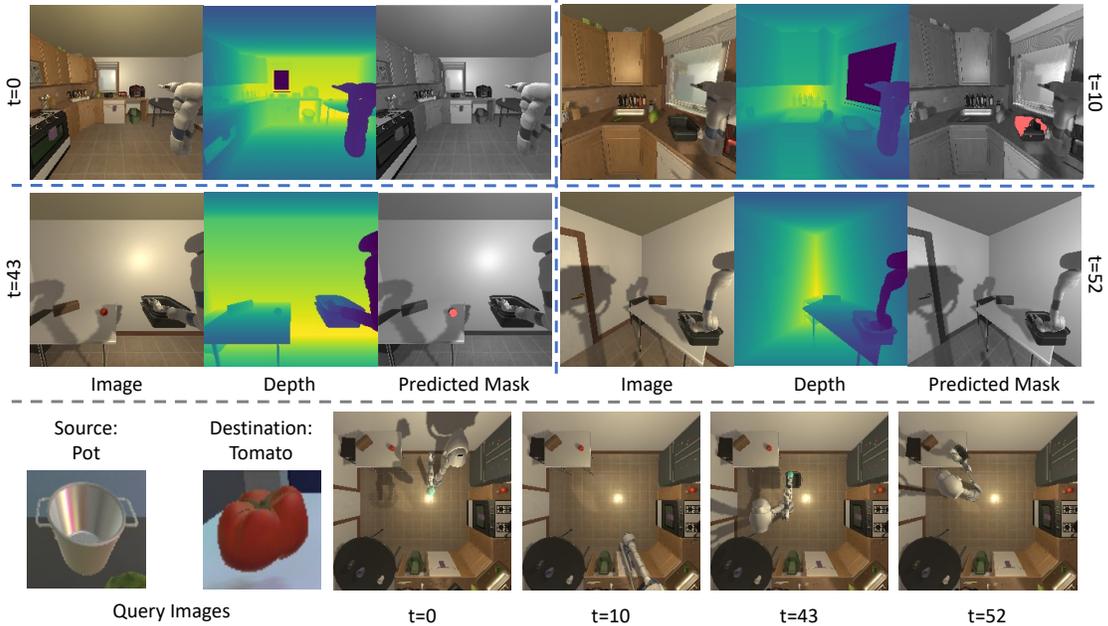
6

Figure 4. **Qualitative Results**. The figure presents a successful episode of *Bringing a Pot to Tomato*. Despite the errors in Pot segmentation and with only a few pixels of the object being segmented, our agent successfully completes the task. Moreover, towards the end of the episode ($t = 52$), the pot in hand occludes the target object (tomato), but the agent is able to remember the target location despite the occlusion. Note that the topdown view (bottom row) is only shown for visualization purposes and is not an input to the network.

compared to our baselines. Rows 5-8 evaluate these models when using ground truth segmentation masks in order to assess these models independent of the segmentation inaccuracies.

Our model m-VOLE, row 4, outperforms all of the baselines (rows 1-3). Directly learning a policy from input images results in a 0 success rate, showing the advantage of using a separate segmentation and target localization network. This result is despite trying various hyperparameters and designs to get this model off the ground. The improvement of the location-aware methods (rows 3 and 4) over the Mask-driven model (row 2) that does not encode relative location demonstrates that the agents benefit from encoding localization information. Our method m-VOLE provides a massive 3x improvement over the baseline in row 3 in terms of Success Rate (SR). This result shows that our approach that estimates the object's distance from the agent is quite effective.

Rows 5-8 evaluate models when employing ground truth segmentation masks. m-VOLE outperforms all baselines significantly, including Loc-MDM with access to the same information as m-VOLE. An interesting observation is that m-VOLE achieves better performance in unseen scenes compared to the ArmPointNav [11] baseline (row 7), which receives the relative location of the target object at each time step. Object Displacement requires attending to the visual observations to avoid collisions of the arm with objects in the scene. Therefore, our conjecture for the higher perfor-

mance is that the ArmPointNav baseline relies heavily on the GPS and location sensors leading to less focus on visual observations. We discuss this further in Appendix D. Figure 4 shows our qualitative results.

One of the main contributions of m-VOLE is the ability to maintain an estimation of the target's location regardless of its visibility. To quantify this contribution, we calculated the percentage of the frames in an episode for which the target is not visible during the inference time. Our experiments show that only in $43.9\%$ and $10.5\%$ of the observed frames the source and goal objects are visible in the agent's frame, respectively (mainly due to arm occlusion or object not being in the frame). More importantly, in $70.5\%$ of those frames, the segmentation model fails to segment the object (misdetection). Despite the low percentage of object visibility, m-VOLE successfully completes the task by aggregating the temporal information.

## 5.2. How robust is m-VOLE to noise?

There are two primary sources of noise that can impact our model, noise in the agent's movements and noise in depth perception. In the following experiments, we evaluate the effect of each one.

**Noise in agent movements.** We use the motion model presented in [32], which models the noise in the traveled distance and angle of rotation. The details of the noise model can be found in Appendix C. Figure 5 shows the robustness of our method to motion noise in comparison
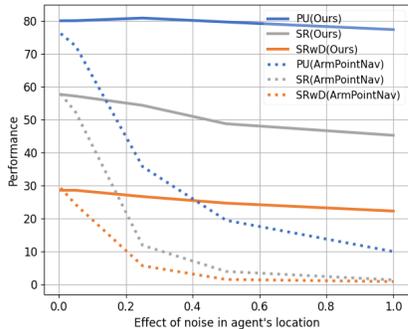
7

Figure 5. **Robustness to noise in agent's movements**. We use the noise model from [32]. The x axis shows the noise multiplier and the y axis shows the performance on the Object Displacement task. Performance without any noise is shown at $x = 0$. $x = 1$ corresponds to a noise as large as the agent step size (20cm).

| Model | PU | SR | SRwD |
|---|---|---|---|
| MDM @100M | 20.5 (-40.0%) | 3.5 (-71.5%) | 2.07 (-72.6%) |
| Loc-MDM @114M | 25.3 (-27.3%) | 4.05 (-64.8%) | 1.89 (-75.1%) |
| m-VOLE (Ours) @20M | **36.3 (-6.2%)** | **6.13 (-47.7%)** | **2.52 (-45.0%)** |

Table 3. **Robustness to noise in depth.** Our model achieves the best performance on all metrics and the lowest relative performance drop compared to the baselines (shown in parentheses).

with ArmPointNav [11] modified to incorporate the same noise model. Our approach m-VOLE is far more robust, whereas ArmPointNav degrades significantly. We believe our method better leverages visual information to recover from noisy estimates.

**Noise in depth perception.** We use the Redwood depth noise model presented in [7]. Table 3 shows the performance of models in the presence of this noise. Note that models are trained with no noise, but inference takes place in the presence of noise. In this experiment, we are not only interested in evaluating models' absolute performance in the presence of noise but also in the relative degradation of each model. With this in mind, to have a fair comparison, we train all models until they reach the same success rate on the training set as m-VOLE. We then evaluate them with and without depth noise and calculate the relative performance drop in all metrics. As seen in Table 3, m-VOLE is more robust to noise in depth compared to our baselines (in absolute terms) and also shows a smaller relative performance degradation. We hypothesize that the aggregated information throughout the episode helps our model recover from the locally observed noise in depth.

## 5.3. Why Conditional Segmentation?

Simultaneously learning to segment objects and learning to plan is challenging. Hence, we isolate the segmentation branch of our network. Most approaches to ObjectNav [4] (navigating towards a specified object) do not use external object detectors and learn to detect and plan jointly, but that is not effective for our task. As we showed in Table 2, the

| Model | PU | SR | SRwD |
|---|---|---|---|
| Our Policy with MaskRCNN [17] detection | 31.3 | 8.7 | 3.87 |
| Our Policy with Conditional Segmentation | **38.7** | **11.6** | **4.59** |

Table 4. **Instance segmentation ablation.** We ablate the performance of our model using different instance segmentation methods. To do so we train a MaskRCNN [17] model on our data and show that our conditional segmentation helps our policy achieve better performance on the ObjDis task.

baseline that does not use a separate segmentation network (row 1) generalizes poorly and achieves 0 success rate.

There are two advantages to our conditional segmentation method compared to the standard state-of-the-art detection/segmentation models such as MaskRCNN [17]. First, our task is more straightforward since we are interested in simply segmenting an instance of the specified object category. Thus, we can afford to train a much smaller network. Second, as shown in Table 4, the model that uses conditional segmentation obtains better performance in all metrics. For this experiment, we use the MaskRCNN [17] model pre-trained on the LVIS dataset [15] and finetuned on AI2-THOR images. We obtain the highest confidence prediction for the target class from MaskRCNN and use that as the input to our policy. This is effective because there is typically just one object instance of the specified category in view.

| Model | Object Set | PU | SR | SRwD |
|---|---|---|---|---|
| Loc-MDM | NovelObj | 9.52 | 0.9 | 0 |
| m-VOLE | NovelObj | **26.2** | **5.24** | **3.33** |
| Loc-MDM | SeenObj | 29.3 | 5.67 | 2.33 |
| m-VOLE | SeenObj | 49.3 | 18.7 | 10.7 |

Table 5. **Zero-shot Manipulation Results.**

## 5.4. Can m-VOLE do zero-shot manipulation?

So far, we have evaluated our model on generalization to unseen scenes. We evaluate whether the model can manipulate novel objects not used for training. This is a challenging task since the variation in object size and shape can result in different types of collisions and occlusions. For this evaluation, we train our model on six object categories (Apple, Bread, Tomato, Lettuce, Pot, Mug) and evaluate on the held-out classes (Potato, Pan, Egg, Spatula, Cup, SoapBottle). Note that our conditional segmentation model that provides the segmentation input to the policy has been trained on all categories. However, our policy has not been exposed to the test objects. Table 5 shows the results. Our model generalizes well to the novel categories, which is challenging as the novel objects' shapes and sizes can result in unseen patterns of collision.

# 6. Conclusion

We propose a method for visual object manipulation, where the goal is to displace an object between two locations in a scene. Our proposed approach learns to estimate target object location via aggregating noisy observations caused by missed detection, view occlusion by the arm, and noisy depth. We show that our approach provides a 3x improvement in success rate over a baseline without this auxiliary information, and it is more robust against noise in depth and agent movements.

# References

[1] Peter Anderson, Angel X. Chang, Devendra Singh Chaplot, Alexey Dosovitskiy, Saurabh Gupta, Vladlen Koltun, Jana Kosecka, Jitendra Malik, Roozbeh Mottaghi, Manolis Savva, and Amir Roshan Zamir. On evaluation of embodied navigation agents. *ArXiv*, 2018. 1, 3

[2] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian D. Reid, Stephen Gould, and Anton van den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *CVPR*, 2018. 1, 3

[3] Dhruv Batra, Angel Xuan Chang, S. Chernova, Andrew J. Davison, Jia Deng, Vladlen Koltun, Sergey Levine, Jitendra Malik, Igor Mordatch, Roozbeh Mottaghi, Manolis Savva, and Hao Su. Rearrangement: A challenge for embodied ai. *arXiv*, 2020. 3

[4] Dhruv Batra, Aaron Gokaslan, Aniruddha Kembhavi, Oleksandr Maksymets, Roozbeh Mottaghi, Manolis Savva, Alexander Toshev, and Erik Wijmans. Objectnav revisited: On evaluation of embodied agents navigating to objects. *arXiv*, 2020. 1, 3, 8

[5] Lillian Y Chang, Siddhartha S Srinivasa, and Nancy S Pollard. Planning pre-grasp manipulation for transport tasks. In *ICRA*, 2010. 2

[6] Devendra Singh Chaplot, Dhiraj Prakashchand Gandhi, Abhinav Gupta, and Russ R Salakhutdinov. Object goal navigation using goal-oriented semantic exploration. In *NeurIPS*, 2020. 1, 3

[7] Sungjoon Choi, Qian-Yi Zhou, and Vladlen Koltun. Robust reconstruction of indoor scenes. In *CVPR*, 2015. 8, 12

[8] Michael Danielczuk, Arsalan Mousavian, Clemens Eppner, and Dieter Fox. Object rearrangement using learned implicit collision functions. In *ICRA*, 2021. 2

[9] Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra. Embodied question answering. In *CVPR*, 2018. 1

[10] Samyak Datta, Oleksandr Maksymets, Judy Hoffman, Stefan Lee, Dhruv Batra, and Devi Parikh. Integrating egocentric localization for more realistic point-goal navigation agents. In *CoRL*, 2020. 3

[11] Kiana Ehsani, Winson Han, Alvaro Herrasti, Eli VanderBilt, Luca Weihs, Eric Kolve, Aniruddha Kembhavi, and Roozbeh Mottaghi. ManipulaTHOR: A Framework for Visual Object Manipulation. In *CVPR*, 2021. 1, 2, 3, 4, 5, 6, 7, 8

[12] Chuang Gan, Jeremy Schwartz, Seth Alter, Martin Schrimpf, James Traer, Julian De Freitas, Jonas Kubilius, Abhishek Bhandwaldar, Nick Haber, Megumi Sano, Kuno Kim, Elias Wang, Damian Mrowca, Michael Lingelbach, Aidan Curtis, Kevin T. Feigelis, Daniel Bear, Dan Gutfreund, David Cox, James J. DiCarlo, Josh H. McDermott, Joshua B. Tenenbaum, and Daniel L. K. Yamins. Threedworld: A platform for interactive multi-modal physical simulation. In *NeurIPS (dataset track)*, 2021. 1

[13] Caelan Reed Garrett, Tomas Lozano-Perez, and Leslie Pack Kaelbling. Ffrob: Leveraging symbolic planning for efficient task and motion planning. *The Intl. J. of Robotics Research*, 2018. 2

[14] Daniel Gordon, Aniruddha Kembhavi, Mohammad Rastegari, Joseph Redmon, Dieter Fox, and Ali Farhadi. Iqa: Visual question answering in interactive environments. In *CVPR*, 2018. 1

[15] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *CVPR*, 2019. 8

[16] Saurabh Gupta, James Davidson, Sergey Levine, Rahul Sukthankar, and Jitendra Malik. Cognitive mapping and planning for visual navigation. In *CVPR*, 2017. 3

[17] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. Mask r-cnn. 2017. 8

[18] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 5

[19] Eric Huang, Zhenzhong Jia, and Matthew T Mason. Large-scale multi-object rearrangement. In *ICRA*, 2019. 2

[20] Yong K Hwang and Narendra Ahuja. Gross motion planning—a survey. *ACM Computing Surveys (CSUR)*, 1992. 2

[21] Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J Davison. Rlbench: The robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 2020. 2

[22] Leslie Pack Kaelbling and Tomás Lozano-Pérez. Integrated task and motion planning in belief space. *The Intl. J. of Robotics Research*, 2013. 2

[23] Jennifer E King, Marco Cognetti, and Siddhartha S Srinivasa. Rearrangement planning using object-centric and robot-centric action spaces. In *ICRA*, 2016. 2

[24] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Daniel Gordon, Yuke Zhu, Abhinav Gupta, and Ali Farhadi. Ai2-thor: An interactive 3d environment for visual ai. *arXiv*, 2017. 1, 3, 5

[25] Athanasios Krontiris and Kostas E Bekris. Dealing with difficult instances of object rearrangement. In *RSS*, 2015. 2

[26] Tomás Lozano-Pérez, Joseph L. Jones, Emmanuel Mazer, and Patrick A. O'Donnell. Task-level planning of pick-and-place robot motions. *Computer*, 1989. 2

[27] Jeffrey Mahler and Ken Goldberg. Learning deep policies for robot bin picking by simulating robust grasping sequences. In *CoRL*, 2017. 2

[28] Pierre Marza, Laëtitia Matignon, Olivier Simonin, and Christian Wolf. Teaching agents how to map: Spatial reasoning for multi-object navigation. *arXiv*, 2021. 3

[29] Piotr W. Mirowski, Razvan Pascanu, Fabio Viola, Hubert Soyer, Andy Ballard, Andrea Banino, Misha Denil, Ross Goroshin, L. Sifre, Koray Kavukcuoglu, Dharshan Kumaran, and Raia Hadsell. Learning to navigate in complex environments. In *ICLR*, 2017. 3

[30] Dipendra Misra, Andrew Bennett, Valts Blukis, Eyvind Niklasson, Max Shatkhin, and Yoav Artzi. Mapping instructions to actions in 3D environments with visual goal prediction. In *EMNLP*, 2018. 3

[31] Tongzhou Mu, Zhan Ling, Fanbo Xiang, Derek Yang, Xuanlin Li, Stone Tao, Zhiao Huang, Zhiwei Jia, and Hao Su. ManiSkill: Generalizable Manipulation Skill Benchmark with Large-Scale Demonstrations. In *NeurIPS (dataset track)*, 2021. 2

[32] Adithyavairavan Murali, Tao Chen, Kalyan Vasudev Alwala, Dhiraj Gandhi, Lerrel Pinto, Saurabh Gupta, and Abhinav Gupta. Pyrobot: An open-source robotics framework for research and benchmarking. *arXiv*, 2019. 7, 8, 12

[33] Srinivas Nedunuri, Sailesh Prabhu, Mark Moll, Swarat Chaudhuri, and Lydia E. Kavraki. Smt-based synthesis of integrated task and motion plans from plan outlines. In *ICRA*, 2014. 2

[34] Matthias Nieuwenhuisen, David Droeschel, Dirk Holz, J. Stückler, Alexander Berner, Jun Yu Li, R. Klein, and Sven Behnke. Mobile bin picking with an anthropomorphic service robot. In *ICRA*, 2013. 2

[35] Santhosh K Ramakrishnan, Ziad Al-Halah, and Kristen Grauman. Occupancy anticipation for efficient exploration and navigation. In *ECCV*, 2020. 3

[36] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large scale visual recognition challenge. *IJCV*, 2015. 5

[37] Nikolay Savinov, Alexey Dosovitskiy, and Vladlen Koltun. Semi-parametric topological memory for navigation. In *ICLR*, 2018. 3

[38] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A platform for embodied ai research. In *ICCV*, 2019. 1, 2

[39] Bokui Shen, Fei Xia, Chengshu Li, Roberto Mart'in-Mart'in, Linxi (Jim) Fan, Guanzhi Wang, S. Buch, Claudia. Pérez D'Arpino, Sanjana Srivastava, Lyne P. Tchapmi, Micael Edmond Tchapmi, Kent Vainio, Li Fei-Fei, and Silvio Savarese. igibson, a simulation environment for interactive tasks in large realistic scenes. In *IROS*, 2021. 3

[40] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P. Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *CVPR*, 2016. 5

[41] Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer,

and Dieter Fox. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In *CVPR*, 2020. 1, 3

[42] Mike Stilman, Jan-Ullrich Schamburek, James Kuffner, and Tamim Asfour. Manipulation planning among movable obstacles. In *ICRA*, 2007. 2

[43] Andrew Szot, Alexander Clegg, Eric Undersander, Erik Wijmans, Yili Zhao, John Turner, Noah Maestre, Mustafa Mukadam, Devendra Singh Chaplot, Oleksandr Maksymets, Aaron Gokaslan, Vladimir Vondrus, Sameer Dharur, Franziska Meier, Wojciech Galuba, Angel Xuan Chang, Zsolt Kira, Vladlen Koltun, Jitendra Malik, Manolis Savva, and Dhruv Batra. Habitat 2.0: Training home assistants to rearrange their habitat. In *NeurIPS*, 2021. 1, 3

[44] Luca Weihs, Matt Deitke, Aniruddha Kembhavi, and Roozbeh Mottaghi. Visual room rearrangement. In *CVPR*, 2021. 1, 3

[45] 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. 5

[46] Erik Wijmans, Abhishek Kadian, Ari S. Morcos, Stefan Lee, Irfan Essa, Devi Parikh, Manolis Savva, and Dhruv Batra. Dd-ppo: Learning near-perfect pointgoal navigators from 2.5 billion frames. In *ICLR*, 2020. 2, 3, 5

[47] Mitchell Wortsman, Kiana Ehsani, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Learning to learn how to learn: Self-adaptive visual navigation using meta-learning. In *CVPR*, 2019. 3

[48] Fei Xia, Chengshu Li, Roberto Martín-Martín, Or Litany, Alexander Toshev, and Silvio Savarese. Relmogen: Leveraging motion generation in reinforcement learning for mobile manipulation. In *ICRA*, 2021. 1, 3

[49] F. Xia, Bokui (William) Shen, Chengshu Li, Priya Kasimbeg, Micael Edmond Tchapmi, Alexander Toshev, Roberto Martín-Martín, and Silvio Savarese. Interactive gibson benchmark: A benchmark for interactive navigation in cluttered environments. *IEEE Robotics and Automation Letters*, 2020. 1, 3

[50] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao Jiang, Yifu Yuan, He Wang, Li Yi, Angel X.Chang, Leonidas Guibas, and Hao Su. SAPIEN: A SimulAted Part-based Interactive ENvironment. In *CVPR*, 2020. 2

[51] Wei Yang, Xiaolong Wang, Ali Farhadi, Abhinav Gupta, and Roozbeh Mottaghi. Visual semantic navigation using scene priors. In *ICLR*, 2019. 3

[52] Lin Yen-Chen, Andy Zeng, Shuran Song, Phillip Isola, and Tsung-Yi Lin. Learning to see before learning to act: Visual pre-training for manipulation. In *ICRA*, 2020. 2

[53] Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine. Metaworld: A benchmark and evaluation for multi-task and meta reinforcement learning. In *CoRL*, 2019. 2

[54] Andy Zeng, Shuran Song, Kuan-Ting Yu, Elliott Donlon, Francois Robert Hogan, Maria Bauzá, Daolin Ma, Orion Taylor, Melody Liu, Eudald Romo, Nima Fazeli, Ferran Alet, Nikhil Chavan Dafle, Rachel Holladay, Isabella

Morona, Prem Qu Nair, Druck Green, Ian Taylor, Weber Liu, Thomas A. Funkhouser, and Alberto Rodriguez. Robotic pick-and-place of novel objects in clutter with multi-affordance grasping and cross-domain image matching. In *ICRA*, 2018. 2

[55] Kuo-Hao Zeng, Luca Weihs, Ali Farhadi, and Roozbeh Mottaghi. Pushing it out of the way: Interactive visual navigation. In *CVPR*, 2021. 3

[56] Xiaoming Zhao, Harsh Agrawal, Dhruv Batra, and Alexander G. Schwing. The surprising effectiveness of visual odometry techniques for embodied pointgoal navigation. *arXiv*, 2021. 3

[57] 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. In *ICRA*, 2017. 1, 2, 3

[58] Yuke Zhu, Josiah Wong, Ajay Mandlekar, and Roberto Martín-Martín. robosuite: A modular simulation framework and benchmark for robot learning. *arXiv*, 2020. 2

## A. Training Details for Conditional Segmentation Model – Sec 4.2

We train our conditional segmentation models for 100 epochs. We use Adam optimizer with a learning rate of $1e - 4$. We initialize the ResNet backbones with ImageNet pre-trained weights for faster training. Input images and query images are of size 224x224. We use random crop and flipping the image for data augmentation during training. We do not use color jittering as we found it to negatively impact the performance of our model on the validation set.

## B. Reward Ablation – Sec 4.3

Section 4.3 introduced three additional reward elements that are not used in the standard ArmPointNav setup. We ablate the performance gain that each component brings to our model. This analysis can be useful not only for evaluating the performance of our method but also for use in future works. Table 6 includes the results in absence of each of these components. Arm control reward, which encourages the agent to move its arm close to the target objects, is shown to be the most important element of the reward and the agent does poorly without this reward. Getting training off the ground with sparse RL reward for object manipulation is extremely difficult, and by motivating the agent to bring its arm closer to the target object, the training becomes faster and more efficient.

| Exploration | Arm Control | Object Visibility | PU | SR | SRwD |
|:---:|:---:|:---:|:---:|:---:|:---:|
| ✓ | - | ✓ | 0.9 | 0 | 0 |
| - | ✓ | ✓ | 47.8 | 15.4 | 5.86 |
| ✓ | ✓ | - | 79.4 | 53.6 | 30.1 |
| ✓ | ✓ | ✓ | **81.2** | **59.6** | **31.0** |

Table 6. **Reward ablations.** We investigate the performance of the model in absence of each reward component. Exploration, arm control, and object visibility refer to the rewards for visiting a new state, $\delta$ distance to object and observing the object, respectively.

## C. Details of the Noise Models – Sec 5.2

This section provides additional details and visualizations on the noise models used in Section 5.2.

### C.1. Agent Motion Noise Model

Murali et al. [32], introduced a noise model for agent's movements and rotation, capturing the noise in motion for a real robot, the Locobot. We use their noise model to ablate the impact of the noise in the agent's movements on our model's performance. To better understand how the noise multiplier (the x-axis in Figure 5 in the main submission)

impacts the predicted trajectory, we illustrate a sample trajectory. Figure 6 shows agent's groundtruth and predicted trajectory for different noise multipliers $x = 0.1 - 1$. Note that the paths start from the same initial location, and the predicted paths diverge from the actual trajectory traversed by the agent as the episode progresses and the errors accumulate.
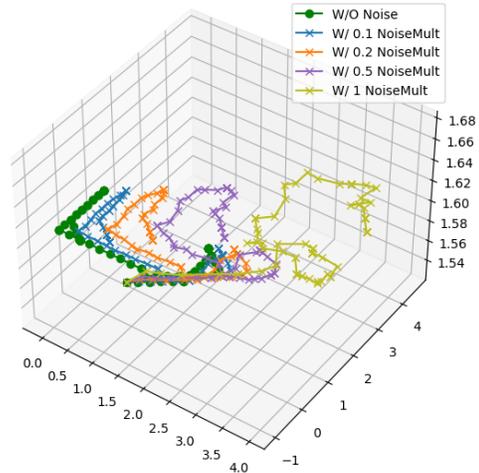


Figure 6. **Noise in Trajectory.** The green path shows the trajectory traversed by the agent, and the other paths show the estimated trajectory of the agent for a variety of noise multipliers. Note that even a small noise multiplier $x = 0.2$ can result in a big divergence in the trajectories.

### C.2. Depth Noise Model

Redwood distortion noise model, introduced in [7], is designed to model the noise of actual depth sensors (Kinect cameras) with a fixed resolution. This noise model can act as a proxy to estimate the impact of using a real-world noisy depth camera.

### C.3. How Does Segmentation Affect the Final Performance? – Sec 5.3

We compare the performance of our approach using different segmentation models in Section 5.3. In this section, we ablate how the performance of the segmentation model affects the final performance of our approach on the task of Object Displacement.

First, we evaluate how the performance on the final task changes if the target objects are detected correctly, but the mask does not perfectly segment the observed object. In other words, if we use IoU as a metric for the accuracy of the segmentation model, how the accuracy of the segmentation network affects the final performance. For this experiment, at each timestep that the object is visible, we randomly choose $x\%$ of the segmentation mask of the target
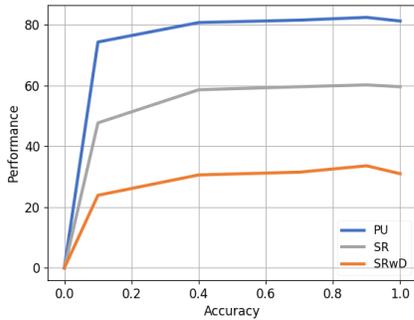
Figure 7. **Impact of Partial Masks.**



Figure 9. **Impact of Category Confusion.**

object to preserve and remove the rest. Figure 7 plots the change in the final performance as we increase $x$. Note that $x = 0$ presents the ablation where no mask is provided, and $x = 1$ is equivalent to providing the groundtruth segmentation mask. This experiment shows that our network can achieve strong results even if only $10\%$ of the target object is segmented.
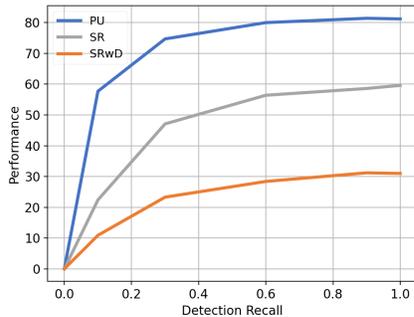


Figure 8. **Impact of Missing Masks.**

We also show that if the segmentation mask of the object is only retrieved $30\%$ of the times, across all the timesteps that the target object is visible, our network can achieve approximately similar performance as the one using the ground-truth segmentation mask (Figure 8). For every $x$ on the plot, at each timestep, we remove the segmentation mask of the target object with the probability $1 - x$ and calculate the final performance of the model. Similar to the previous plot, $x = 0$ shows the performance when no mask is provided, and $x = 1$ is equivalent to the groundtruth segmentation mask.

One of the other main issues that the segmentation models have is detecting wrong objects. For instance, as shown in the supplementary video, the segmentation model might segment a pan in the scene while the query object asks for a pot. Figure 9 shows the final performance of the model for different rates of mis-detections. At each timestep, with the
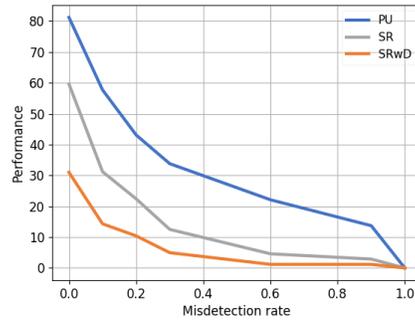
probability of $x\%$, instead of the segmentation mask of the target object, we randomly select the segmentation mask of another object in the scene as the input mask to the model. For instance, metrics at $x = 0.2$ show the model's performance using a segmentation network that detects a wrong object with the probability of $20\%$.

The conclusion is that even segmentation models with high precision and low recall are beneficial for the Object Displacement task.

## D. Our method (m-VOLE) outperforming ArmPointNav (APN). – Sec 5.1



Figure 10. **m-VOLE vs APN paths towards the target.** The red path illustrates the greedy solution an agent with access to GT direction sensors might take to reach the bowl (Failing due to collision with the oven).

APN has access to the exact direction towards the target, so it tends to choose the direct path towards the goal, which is not necessarily a plausible solution. For example, as shown in Fig. 10, the red path represents a shorter path, but the arm collides with other objects if it follows that path. Moreover, we did a further quantitative analysis to investigate m-VOLE superiority to APN (Tab 7). EpLen PU and EpLen Success, respectively, represent episode length for the pickup stage and the full task. Our analysis shows that these metrics are lower for APN than m-VOLE. However, the SRwD of APN (success without collision) is lower. This

13

observation shows that APN is more efficient in the number of steps (as it uses the exact direction) but does not handle collisions well.

| Model | EpLen PU | EpLen Success | SRwD |
|---|---|---|---|
| ArmPointNav | 46.3 | 78.8 | 28.5 |
| m-VOLE w/ GT mask | 49.3 | 87.1 | 31.0 |

Table 7. **m-VOLE vs APN.**

## E. Limitations

Here, we discuss two main limitations of the work. First, we consider separate modules for segmentation and depth perception. Hence, their errors have a cascading effect. Another design choice we made was to abstract away grasping. While this allowed us to focus on other challenging aspects of the problem, moving beyond simulation will likely require methods to account for real-world graspers. These are some of our study's primary limitations, which serve as avenues for future work.