

# OK-Robot: What Really Matters in Integrating Open-Knowledge Models for Robotics

Peiqi Liu\*<sup>1</sup>

Yaswanth Orru\*<sup>1</sup>

Chris Paxton<sup>2</sup>

Nur Muhammad Mahi Shafiullah<sup>†1</sup>

Lerrel Pinto<sup>†1</sup>

New York University<sup>1</sup>, AI at Meta<sup>2</sup>

<https://ok-robot.github.io>

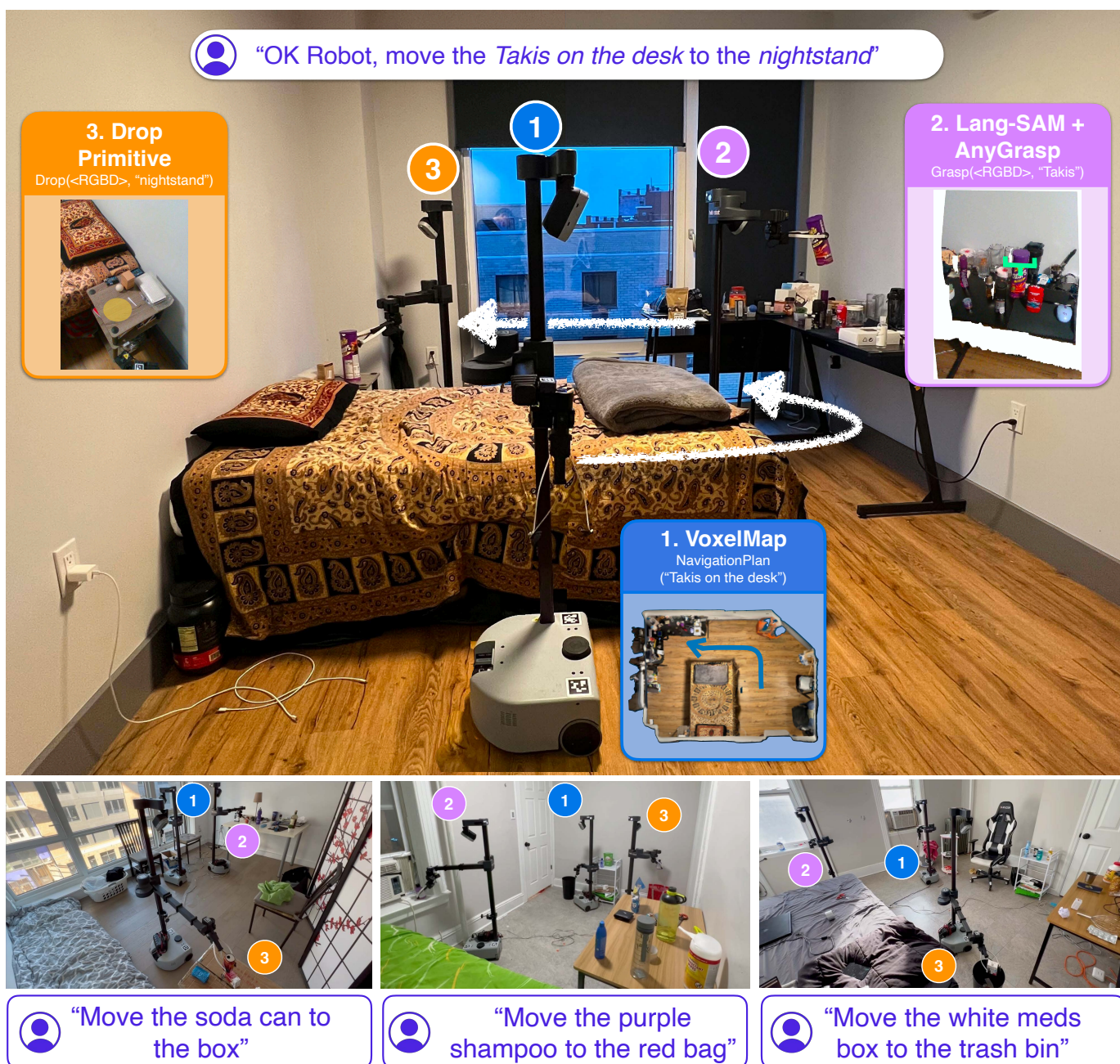


Fig. 1: OK-Robot is an Open Knowledge robotic system, which integrates a variety of learned models trained on publicly available data, to pick and drop objects in real-world environments. Using Open Knowledge models such as CLIP, Lang-SAM, AnyGrasp, and OWL-ViT, OK-Robot achieves a 58.5% success rate across 10 unseen, cluttered home environments, and 82.4% on cleaner, decluttered environments.

**Abstract**— Remarkable progress has been made in recent years in the fields of vision, language, and robotics. We now have vision models capable of recognizing objects based on language queries, navigation systems that can effectively control mobile systems, and grasping models that can handle a wide range of objects. Despite these advancements, general-purpose applications of robotics still lag behind, even though they rely on these fundamental capabilities of recognition, navigation, and grasping. In this paper, we adopt a systems-first approach to develop a new Open Knowledge-based robotics framework called OK-Robot. By combining Vision-Language Models (VLMs) for object detection, navigation primitives for movement, and grasping primitives for object manipulation, OK-Robot offers an integrated solution for pick-and-drop operations without requiring any training. To evaluate its performance, we run OK-Robot in 10 real-world home environments. The results demonstrate that OK-Robot achieves a 58.5% success rate in open-ended pick-and-drop tasks, representing a new state-of-the-art in Open Vocabulary Mobile Manipulation (OVMM) with nearly  $1.8\times$  the performance of prior work. On cleaner, uncluttered environments, OK-Robot’s performance increases to 82%. However, the most important insight gained from OK-Robot is the critical role of nuanced details when combining Open Knowledge systems like VLMs with robotic modules.

## I. INTRODUCTION

Creating a general-purpose robot has been a longstanding dream of the robotics community. With the increase in data-driven approaches and large robot models, impressive progress is being made [1–4]. However, current systems are brittle, closed, and fail when encountering unseen scenarios. Even the largest robotics models can often only be deployed in previously seen environments [5, 6]. The brittleness of these systems is further exacerbated in settings where little robotic data is available, such as in unstructured home environments.

The poor generalization of robotic systems lies in stark contrast to large vision models [7–10], which show capabilities of semantic understanding [11–13], detection [7, 8], and connecting visual representations to language [9, 10, 14]. At the same time, base robotic skills for navigation [15], grasping [16–19], and rearrangement [20, 21] are fairly mature. Hence, it is perplexing that robotic systems that combine modern vision models with robot-specific primitives perform so poorly. To highlight the difficulty of this problem, the recent NeurIPS 2023 challenge for open-vocabulary mobile manipulation (OVMM) [22] registered a success rate of 33% for the winning solution [23].

So what makes open-vocabulary robotics so hard? Unfortunately, there isn’t a single challenge that makes this problem hard. Instead, inaccuracies in different components compound and together results in overall drop. For example, the quality of open-vocabulary retrievals of objects in homes is dependent on the quality of query strings, navigation targets determined from VLMs may not be reachable to the robot, and the choice of different grasping models may lead to large differences in grasping performance. Hence, making progress on this problem requires a careful and nuanced framework that both integrates

VLMs and robotics primitives, while being flexible enough to incorporate newer models as they are developed by the VLM and robotics community.

We present OK-Robot, an *Open Knowledge Robot* that integrates state-of-the-art VLMs with powerful robotics primitives for navigation and grasping to enable pick-and-drop. Here, *Open Knowledge* refers to learned models trained on large, publicly available datasets. When placed in a new home environment, OK-Robot is seeded with a scan taken from an iPhone. Given this scan, dense vision-language representations are computed using LangSam [24] and CLIP [9] and stored in a semantic memory. Then, given a language-query for an object that has to be picked, language representations of the query is matched with semantic memory. After this, navigation and picking primitives are applied sequentially to move to the desired object and pick it up. A similar process can be carried out for dropping the object.

To study OK-Robot, we tested it in 10 real world home environments. Through our experiments, we found that on a never seen, natural home environment, a zero-shot deployment of our system achieves 58.5% success on average. However, this success rate is largely dependant on the “naturalness” of the environment, as we show that with improving the queries, decluttering the space, and excluding objects that are clearly adversarial (too large, too translucent, too slippery), this success rate reaches about 82.4%. Overall, through our experiments, we make the following observations:

- **Pre-trained VLMs are highly effective for open-vocabulary navigation:** Current open-vocabulary vision-language models such as CLIP [9] or OWL-ViT [8] offer strong performance in identifying arbitrary objects in the real world, and enable navigating to them in a zero-shot manner (see Section II-A.)
- **Pre-trained grasping models can be directly applied to mobile manipulation:** Similar to VLMs, special purpose robot models pre-trained on large amounts of data can be applied out of the box to approach open-vocabulary grasping in homes. These robot models do not require any additional training or fine-tuning (see Section II-B.)
- **How components are combined is crucial:** Given the pretrained models, we find that they can be combined with no training using a simple state-machine model. We also find that using heuristics to counteract the robot’s physical limitations can lead to a better success rate in the real world (see Section II-D.)
- **Several challenges still remain:** While, given the immense challenge of operating zero-shot in arbitrary homes, OK-Robot improves upon prior work, by analyzing the failure modes we find that there are significant improvements that can be made on the VLMs, robot models, and robot morphology, that will directly increase performance of open-knowledge manipulation agents (see Section III-D.)

To encourage and support future work in open-knowledge robotics, we will share the code and modules for OK-Robot, and are committed to supporting reproduction of our results.

\* Denotes equal contribution and † denotes equal advising.  
Correspondence to: mahi@cs.nyu.edu

More information along with robot videos are available on our project website: <https://ok-robot.github.io>.

## II. TECHNICAL COMPONENTS AND METHOD

Our method, on a high level, solves the problem described by the query: “Pick up **A** (from **B**) and drop it on/in **C**”, where **A** is an object and **B** and **C** are places in a real-world environment such as homes. The system we introduce is a combination of three primary subsystems combined on a Hello Robot: Stretch. Namely, these are the open-vocabulary object navigation module, the open-vocabulary RGB-D grasping module, and the dropping heuristic. In this section, we describe each of these components in more details.

### A. Open-home, open-vocabulary object navigation

The first component of our method is an open-home, open-vocabulary object navigation model that we use to map a home and subsequently navigate to any object of interest designated by a natural language query.

**Scanning the home:** For open vocabulary object navigation, we follow the approach from CLIP-Fields [27] and assume a pre-mapping phase where the home is “scanned” manually using an iPhone. This manual scan simply consists of taking a video of the home using the Record3D app on the iPhone, which results in a sequence of posed RGB-D images.

Alternatively, this could be done automatically using frontier-based exploration [15, 25, 26], but for speed and simplicity we prefer the manual approach [26, 27]. We take this approach since the frontier-based approaches tend to be slow and cumbersome, especially for a novel space, while our “scan” take less than one minute for each room. Once collected, the RGB-D images, along with the camera pose and positions, are exported to our library for map-building.

To ensure our semantic memory contains both the objects of interest as well as the navigable surface and any obstacles, the recording must capture the floor surface alongside the objects and receptacles in the environment.

**Detecting objects:** On each frame of the scan, we run an open-vocabulary object detector. Unlike previous works which used Detic [7], we chose OWL-ViT [8] as the object detector since we found it to perform better in preliminary queries. We apply the detector on every frame, and extract each of the object bounding box, CLIP-embedding, detector confidence, and pass them onto the object memory module of our navigation module.

Building on top of previous work [27], we further refine the bounding boxes into object masks with Segment Anything (SAM) [28]. Note that, in many cases, open-vocabulary object detectors still require a set of natural language object queries that they try to detect. We supply a large set of such object queries, derived from the original Scannet200 labels [29], such that the detector captures most common objects in the scene.

**Object-centric semantic memory:** We use an object-centric memory similar to Clip-Fields [27] and OVMM [25] that we call the VoxelMap. The object masks are back-projected in real-world coordinates using the depth image and the pose

collected by the camera, giving us a point cloud where each point has an associated semantic vector coming from CLIP. Then, we voxelize the point cloud to a 5 cm resolution and for each voxel, calculate the detector-confidence weighted average for the CLIP embeddings that belong to that voxel. This voxel map builds the base of our object memory module. Note that the representation created this way remains static after the first scan, and cannot be adapted during the robot’s operation. This inability to dynamically create a map is discussed in our limitations section (Section V).

**Querying the memory module:** Our semantic object memory gives us a static world representation represented as possibly non-empty voxels in the world, and a semantic CLIP vector associated with each voxel. Given a language query, we convert it to a semantic vector using the CLIP language encoder. Then, we find the top voxel where the dot product between the encoded vector and the voxel’s semantic representation is maximized. Since each voxel is associated with a real location in the home, this lets us find the location where a queried object is most likely to be found, similar to Figure 2(a).

When necessary, we implement “A on B” as “A near B”. We do so by selecting top-10 points for query A and top-50 points for query B. Then, we calculate the  $10 \times 50$  pairwise euclidean distances and pick the A-point associated with the shortest (A, B) distance. Note that, during the object navigation phase, we use this query only to navigate to the object approximately, and not for manipulation.

This approach gives us two advantages: our map can be at lower resolution than those in prior work [26, 27, 30], and we can deal with small movements in object’s location after building the map.

**Navigating to objects in the real world:** Once our navigation model gives us a 3D location coordinate in the real world, we use that as a navigation target for our robot to initialize our manipulation phase. In previous works [15, 27, 31], the navigation objective was to go and look at an object, which can be done while staying at a safe distance from the object itself. In contrast, our navigation module must place the robot at an arms length so that the robot can manipulate the target object afterwards. Thus, our navigation method has to balance the following objectives:

- 1) the robot needs to be close enough to the object to manipulate it,
- 2) the robot needs some space to move its gripper, so there needs to be a small but non-negligible space between the robot and the object and
- 3) the robot needs to avoid collision during manipulation, and thus needs to keep its distance from all obstacles.

We use three different navigation score functions, each associated with one of our previous concerns, and evaluate them on each point of the space to find the best position to place the robot.

Let a random point be  $\vec{x}$ , the closest obstacle point as  $\vec{x}_{obs}$ , and the target object as  $\vec{x}_o$ . Then, we can define the following three functions  $s_1, s_2, s_3$  to capture our three criterion. Then,

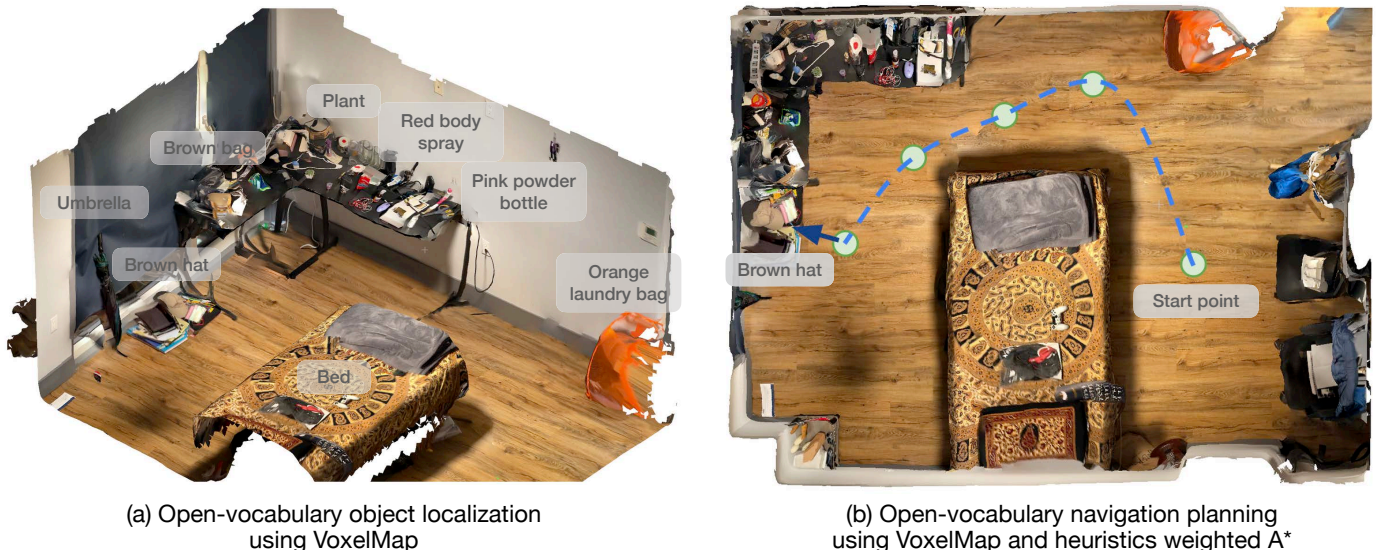


Fig. 2: Open-vocabulary, open knowledge object localization and navigation in the real-world. We use the VoxelMap [25] for localizing objects with natural language queries, and use an A\* algorithm similar to USANet [26] for path planning.

their weighted sum  $s$  to find the ideal navigation point  $\vec{x}^*$  in our space that minimizes  $s(\vec{x})$ , and the direction is towards the vector from  $\vec{x}^*$  to  $\vec{x}_o$ .

$$\begin{aligned}
 s_1(\vec{x}) &= \|\vec{x} - \vec{x}_o\| \\
 s_2(\vec{x}) &= 40 - \min(\|\vec{x} - \vec{x}_o\|, 40) \\
 s_3(\vec{x}) &= \begin{cases} 1/\|\vec{x} - \vec{x}_{obs}\|, & \text{if } \|\vec{x} - \vec{x}_{obs}\|_0 \leq 30 \\ 0, & \text{otherwise} \end{cases} \\
 s(\vec{x}) &= s_1(\vec{x}) + 8s_2(\vec{x}) + 8s_3(\vec{x})
 \end{aligned}$$

To navigate to this target point safely from any other point in space, we follow a similar approach to [26, 32] by building an obstacle map from our previously captured posed RGB-D images. We build a 2D, 10cm×10cm grid of obstacles over which we navigate using the A\* algorithm. To convert our voxel map to an obstacle map, we first set a floor and ceiling height. Presence of occupied voxels in between them implies the grid cell is occupied, while presence of neither ceiling nor floor voxels mean that the grid cell is unexplored. We mark both occupied or unexplored cells as not navigable. Around each occupied point, we mark any point within a 20 cm radius as also non-navigable to account for the robot’s radius and a turn radius. In our A\* algorithm, we use the  $s_3$  function as a heuristic on the node costs to navigate further away from any obstacles, which makes our generated paths similar to ideal Voronoi paths [33] in our experiments.

*B. Open-vocabulary grasping in the real world*

Unlike open-vocabulary navigation, for grasping, our method needs to physically interact with arbitrary objects in the real world, which makes this part significantly more difficult. As a result, we opt for using a pre-trained grasping model to generate grasp poses in the real world, and filter that with language-conditioning using a modern VLM.

**Grasp perception:** Once the robot reaches the object using the navigation method outlined in Section II-A, we use a pre-trained grasping model or heuristic to generate a grasp for the robot. We point the robot’s RGB-D head camera towards the object’s location in space, as given to us by our semantic memory module, and capture an RGB-D image from it (Figure 3, column 1). We backproject and convert the depth image to a pointcloud as necessary. Then, we pass this information to our grasp generation module. The grasp generation module that we use in our work is AnyGrasp [19], which generates collision free grasps with a parallel jaw gripper in a scene given a single RGB image and a pointcloud.

AnyGrasp provides us with possible grasps in the scene (Figure 3 column 2) with grasp point, width, height, depth, and a “graspness score”, which indicates uncalibrated model confidence in each grasp. However, such modules generally generate all possible grasps in a scene, which we need to filter using the language query.

**Filtering grasps using language queries:** Once we get all proposed grasps from AnyGrasp, we filter the grasps using LangSam [24]. We use LangSam [24] to segment the captured image and get the desired object’s mask with the language query (Figure 3 column 3). Then, we project all the proposed grasp points onto the image and find the grasps that fall into the object mask (Figure 3 column 4).

We pick the best grasp using a heuristic, where if the grasp score is  $\mathcal{S}$  and the angle between the grasp normal and floor normal is  $\theta$ , then the new heuristic score is  $\mathcal{S} - (\theta^4/10)$ . This heuristic prioritizes grasps with the highest graspness score but also a horizontally flat proposed grasp. We prefer horizontal grasps because they are robust to small calibration errors on the robot, while vertical grasps need to be quite point-accurate to be successful. Being robust to hand-eye calibration errors is a desired property as we transport the robot to different homes

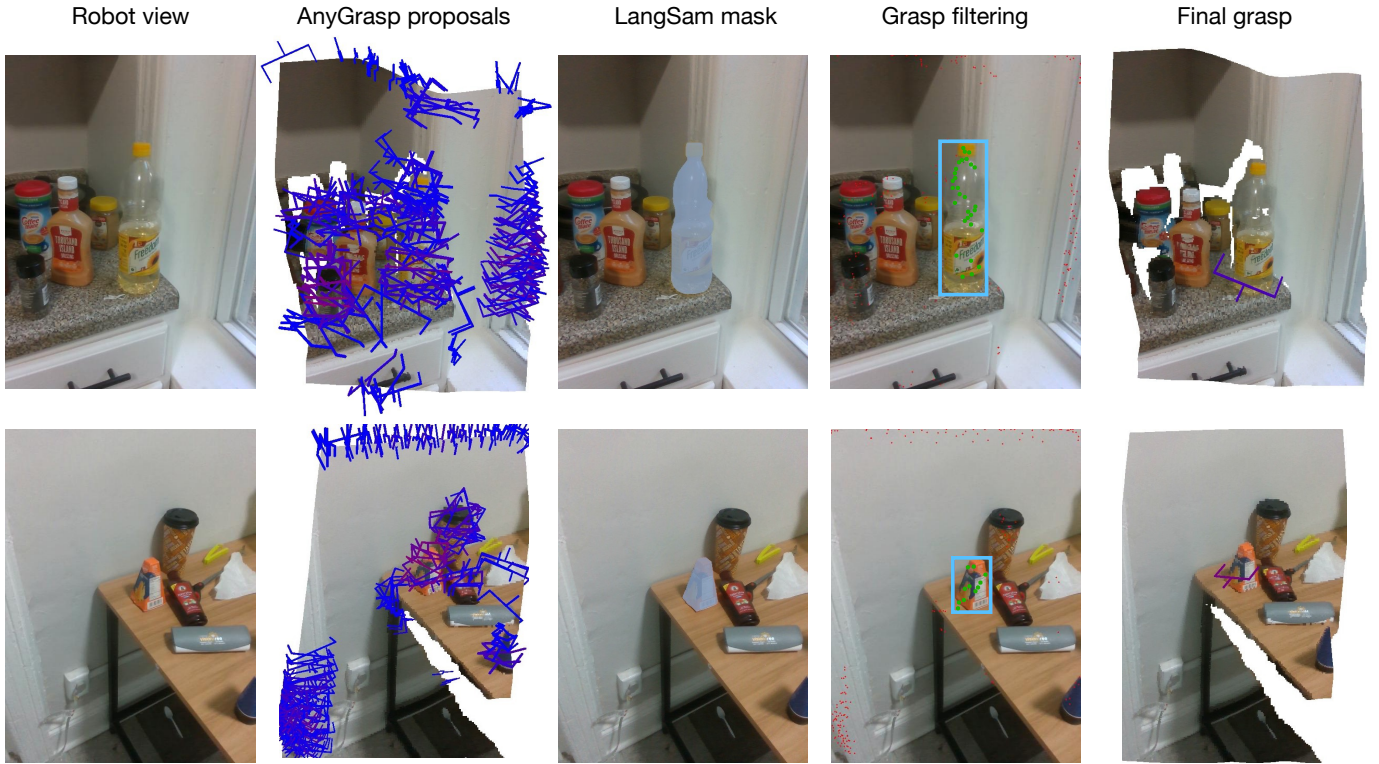


Fig. 3: Open-vocabulary grasping in the real world. From left to right, we show the (a) robot POV image, (b) all suggested grasps from AnyGrasp [19], (c) object mask given label from LangSam [24], (d) grasp points filtered by the mask, and (e) grasp chosen for execution.

over the course of our experiments.

**Grasp execution:** Once we identify the best grasp (Figure 3 column 5), we use a simple pre-grasp approach [34] to grasp our intended object. Let’s assume that  $\vec{p}$  is the grasp point and  $\vec{a}$  is the approach vector given by our grasping model. Then, our robot gripper follows the following trajectory:

$$\langle \vec{p} - 0.2\vec{a}, \vec{p} - 0.08\vec{a}, \vec{p} - 0.04\vec{a}, \vec{p} \rangle$$

Put simply, our method approaches the object from a pre-grasp position in a line with progressively smaller motions. Moving slower as we approach the object is important since the robot can knock over light objects otherwise. Once we reach the predicted grasp point, we close the gripper in a close loop fashion to make sure we can get a solid grip on the object without crushing it. Finally, after grasping the object, we lift up the robot arm, retract it fully, and rotate the wrist to have the object tucked over the body. This behavior maintains the robot footprint while ensuring the object is held securely by the robot and doesn’t fall while navigating to the drop location.

### C. Dropping heuristic

After picking up an object, we find and navigate to the location to drop it using the same methods as described in Section II-A. Unlike in HomeRobot’s baseline implementation [25], which assumes that the drop-off location is a flat surface, we extend our heuristic to also cover concave objects such as sink, bins, boxes, and bags. First, we segment the point cloud  $P$  captured by the robot’s head camera using

LangSam [24] similar to Section II-B using the drop language query. Then, we align that segmented point cloud such that X-axis is aligned with the way the robot is facing, Y-axis is to its left and right, and the Z-axis of the point cloud is aligned with the floor normal. We call this aligned pointcloud  $P_a$ . Finally, we normalize the point cloud so that the robot’s X- and Y- coordinate is (0, 0), and the floor plane is at  $z = 0$ . On the aligned, segmented point cloud, we consider the X- and Y-coordinates for each point, and find the respective medians on each axis that we call  $x_m$  and  $y_m$ . Finally, we find a drop height using  $z_{\max} = 0.2 + \max\{z \mid (x, y, z) \in P_a; 0 \leq x \leq x_m; |y - y_m| < 0.1\}$  on the segmented, aligned pointcloud. We add a small buffer of 0.2 to the height to avoid collisions between the robot and the drop location. Finally, we move the robot gripper above the drop point, and open the gripper to drop the object. While this heuristic sometimes fails to place an object on a cluttered surface, in our experiments it performs well on average.

### D. Deployment in homes

Once we have our navigation, pick, and drop primitive in place, we combine them directly to create our robot method that can be applied in any novel home directly. For a new home environment, we can “scan” the room in under a minute. Then, it takes less than five minutes to process that into our VoxelMap. For our ablations, it takes about 50 minutes to train the necessary implicit semantic fields/SDF models such as CLIP-Fields or USA-Net if we are using them. Once that is

done, the robot can be immediately placed at the base and start operating. From arriving into a completely novel environment to start operating autonomously in it, our system takes under 10 minutes on average to complete the first pick-and-drop task.

**State machine model:** The transition between different modules happens automatically, in a predefined fashion, once a user specifies the object to pick and where to drop it. Since we do not implement error detection or correction, our state machine model is a simple linear chain of steps leading from navigating to object, to grasping, to navigating to goal, and to dropping the object at the goal to finish the task.

**Protocol for home experiments:** To run our experiment in a novel home, we first move the robot to a previously unobserved room. There, we record the scene and create our VoxelMap. Concurrently, we arbitrarily pick between 10-20 objects in each scene that can fit in the robot gripper. These are objects “found” in the scene, and are not ones selected beforehand. We come up with a language query for each chosen object using GPT-4V [35] to keep the queries consistent and free of experimenter bias. The effect of different queries for the same object on OK-Robot is discussed in Section III-D. Then, we query our navigation module to filter out all the navigation failures; i.e. objects whose location could not be found by our semantic memory module. Then, we execute pick-and-drop on remaining objects sequentially, without resets between trials.

### III. EXPERIMENTS

We evaluate our method in two set of experiments. On the first set of experiments, we evaluate between multiple alternatives for each of our navigation and manipulation modules. These experiments give us insights about which modules to use and evaluate in a home environment as a part of our method. On the next set of experiments, we took our robots to 10 homes and ran 171 pick-and-drop experiments to empirically evaluate how our method performs in completely novel homes, and to understand the failure modes of our system.

Through these experiments, we look to answer a series of questions regarding the capabilities and limits of current Open Knowledge robotic systems, as embodied by OK-Robot. Namely, we ask the following:

- 1) How well can such a system tackle the challenge of pick and drop in arbitrary homes?
- 2) How well do alternate primitives for navigation and grasping compare to the recipe presented here for building an Open Knowledge robotic system?
- 3) How well can our current systems handle unique challenges that make homes particularly difficult, such as clutter, ambiguity, and affordance challenges?
- 4) What are the failure modes of such a system and its individual components in real home environments?

#### A. List of home experiments

Over the 10 home environment, OK-Robot achieved a 58.5% success rates in completing full pick-and-drops. Notably, this is a zero-shot algorithm, and the success rate is over novel

objects sourced from each home. As a result, each of the success and the failure of the robot tells us something interesting about applying open-knowledge models in robotics, which is what we analyze over the next sections.

In Appendix C, we provide the details of all our home experiments and results from the same, and in Appendix B we show a subset of the objects OK-Robot operated on. Snippets of our experiments are in Figure 1, and full videos can be seen on our project website.

#### B. Ablations over system components

Apart from the navigation and manipulation strategies that we used in the home experiments, we also evaluated a number of alternative semantic memory module and open vocabulary navigation modules. We compared them by evaluating them in three different environment setups in our lab.

**Alternate semantic navigation strategies:** We evaluate the following semantic memory modules:

- **VoxelMap [25]:** VoxelMap converts every detected object to a semantic vector and stores such info into an associated voxel. Occupied voxels serve as an obstacle map.
- **CLIP-Fields [27]:** CLIP-Fields converts a sequence of posed RGB-D images to a semantic vector field by using open-label object detectors and semantic language embedding models. The result associates each point in the space with two semantic vectors, one generated via a VLM [9], and another generated via a language model [36], which is then embedded into a neural field [37].
- **USA-Net [26]:** USA-Net generates multi-scale CLIP features and embeds them in a neural field that also doubles as a signed distance field. As a result, a single model can support both object retrieval and navigation.

We compare them in the same three environments with a fixed set of queries, the results of which are shown in Figure 5.

**Alternate grasping strategies:** Similarly, we compare multiple grasping strategies to find out the best grasping strategy for our method.

- **AnyGrasp [19]:** AnyGrasp is a single view RGB-D based grasping model. It is trained on the GraspNet dataset which contains 1B grasp labels.
- **Open Graspness [19]:** Since the AnyGrasp model is free but not open source, we use an open licensed baseline trained on the same dataset.
- **Contact-GraspNet [16]:** We use Contact-GraspNet as a prior work baseline, which is trained on the Acronym [38] dataset. One limitation of Contact-GraspNet is that it was trained on a fixed camera view for a tabletop setting. As a result, in our application with a moving camera and arbitrary locations, it failed to give us meaningful grasps.
- **Top-down grasp [25]:** As a heuristic based baseline, we compare with the top-down heuristic grasp provided in the HomeRobot project.

In Figure 5, we see their comparative performance in three lab environments. For semantic memory modules, we see that VoxelMap, used in OK-Robot and described in Sec. II-A,

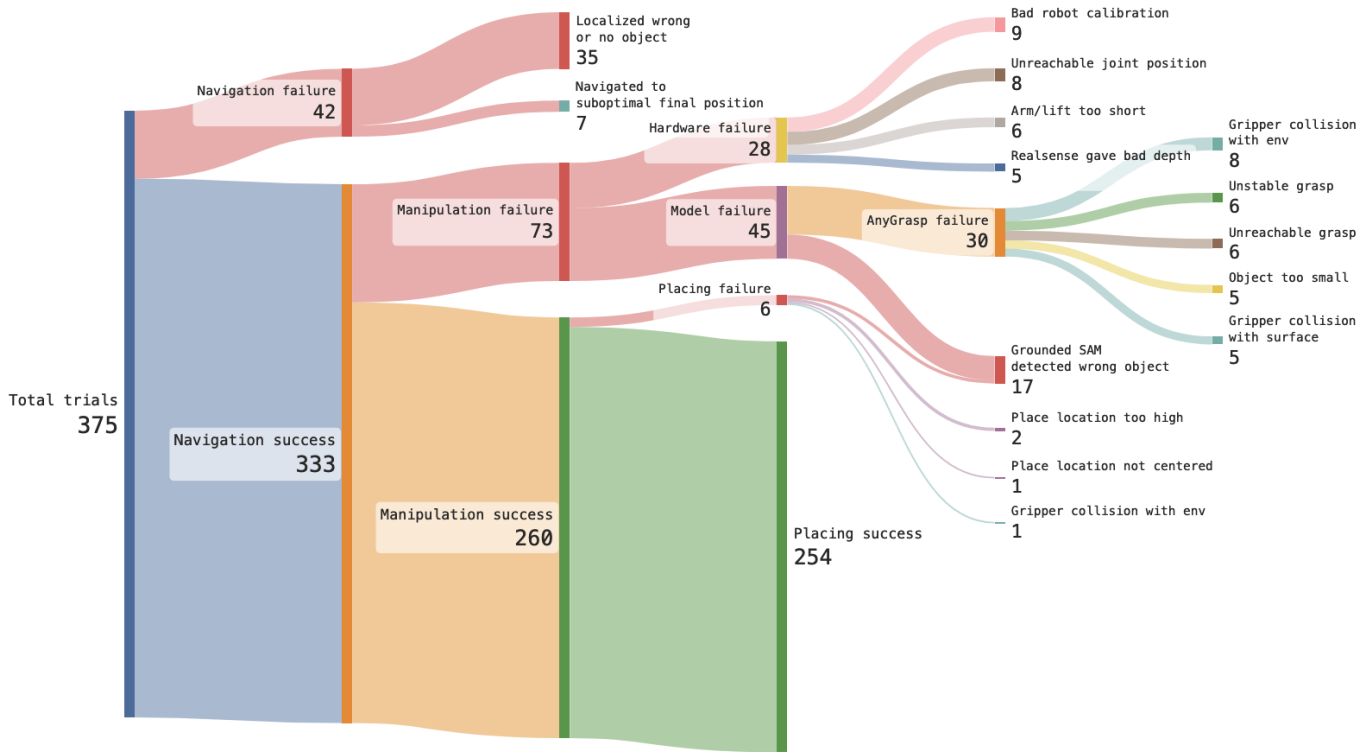


Fig. 4: All the success and failure cases in our home experiments, aggregated over all three cleaning phases, and broken down by mode of failure. From left to right, we show the application of the three components of OK-Robot, and show a breakdown of the long-tail failure modes of each of the components.

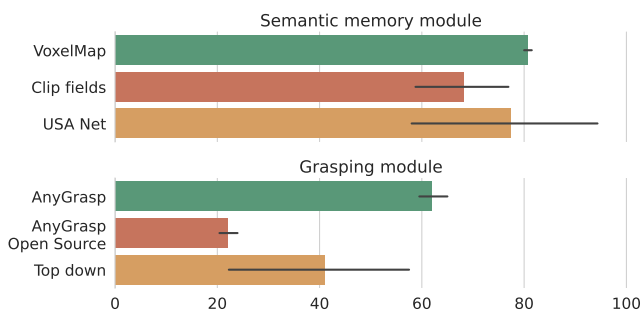


Fig. 5: Ablation experiment using different semantic memory and grasping modules, with the bars showing average performance and the error bars showing standard deviation over the environments.

outperforms other semantic memory modules by a small margin. It also has much lower variance compared to the alternatives, meaning it is more reliable. As for grasping modules, AnyGrasp clearly outperforms other grasping methods, performing almost 50% better in a relative scale over the next best candidate, top-down grasp. However, the fact that a heuristic-based algorithm, top-down grasp from HomeRobot [25] beats the open-source AnyGrasp baseline and Contact-GraspNet shows that building a truly general-purpose grasping model remains difficult.

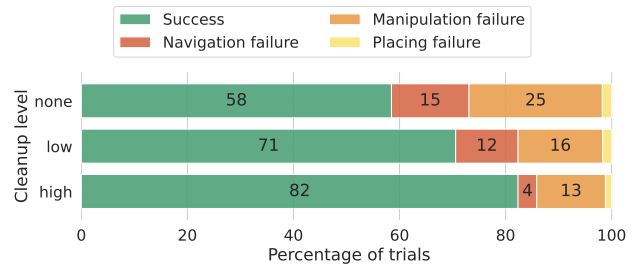


Fig. 6: Failure modes of our method in novel homes, broken down by the failures of the three modules and the cleanup levels.

### C. Impact of clutter, object ambiguity, and affordance

What makes home environments especially difficult compared to lab experiments is the presence of physical clutter, language-to-object mapping ambiguity, and hard-to-reach positions. To gain a clear understanding of how such factors play into our experiments, we go through two “clean-up” processes in each environment. During the clean-up, we pick a subset of objects that are free from ambiguity from the previous rounds, clean the clutter around objects, and generally relocated them in an accessible locations. We go through two of such clean-up rounds at each environment, which gives us insights about the performance gap caused by the natural difficulties of a home-like environment.

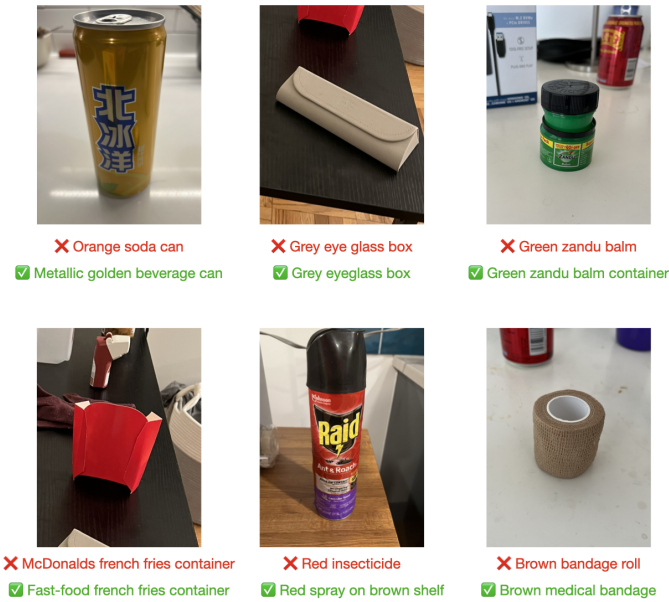


Fig. 7: Samples of failed or ambiguous language queries into our semantic memory module. Since the memory module depends on pretrained large vision language model, its performance shows susceptibility to particular “incantations” similar to current LLMs.

We show a complete analysis of the tasks listed section III-A which failed in various stages in Figure 6. As we can see from this breakdown, as we clean up the environment and remove the ambiguous objects, the navigation accuracy goes up, and the total error rate goes down from 15% to 12% and finally all the way down to 4%. Similarly, as we clean up clutter from the environment, we find that the manipulation accuracy also improves and the error rates decrease from 25% to 16% and finally 13%. Finally, since the drop-module is agnostic of the label ambiguity or manipulation difficulty arising from clutter, the failure rate of the dropping primitive stays roughly constant through the three phases of cleanup.

#### D. Understanding the performance of OK-Robot

While our method can show zero-shot generalization in completely new environments, we probe OK-Robot to better understand its failure modes. Primarily, we elaborate on how our model performed in novel homes, what were the biggest challenges, and discuss potential solutions to them.

We first show a coarse-level breakdown of the failures, only considering the three high level modules of our method in Figure 6. We see that generally, the leading cause of failure is our manipulation failure, which intuitively is the most difficult as well. However, at a closer look, we notice a long tail of failure causes, which is presented in figure 4.

We see that the leading three cause of failures are failing to retrieve the right object to navigate to from the semantic memory (9.3%), getting a difficult pose from the manipulation module (8.0%), and hardware difficulties (7.5%). In this section, we go over the analysis of the failure modes presented in Figure 4 and discuss the most frequent cases.

**Natural language queries for objects:** One of the primary

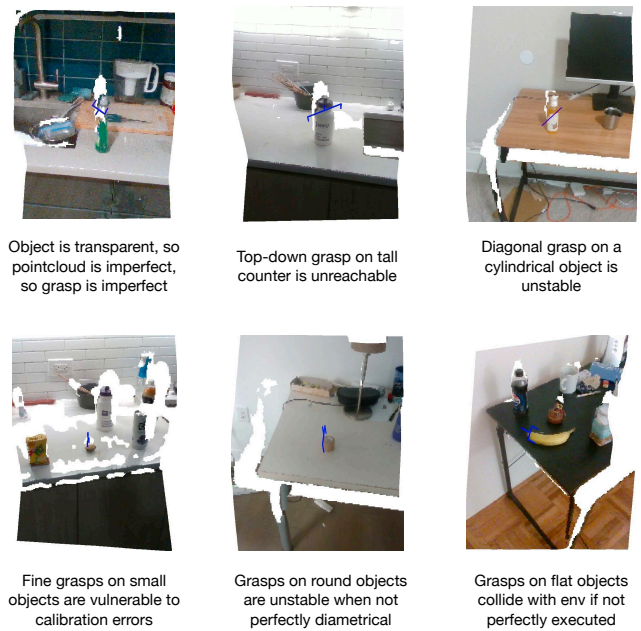


Fig. 8: Samples of failures of our manipulation module. Most failures stem from using only a single RGB-D view to generate the grasp and the limiting form-factor of a large two-fingered parallel jaw gripper.

reasons our OK-Robot can fail is when a natural language query given by the user doesn’t retrieve the intended object from the semantic memory. In Figure 7 we show how some queries may fail while semantically very similar but slightly modified wording of the same query might succeed.

Generally, this has been the case for scenes where there are multiple visually or semantically similar objects, as shown in the figure. There are other cases where some queries may pass while other very similar queries may fail. An interactive system that gets confirmation from the user as it retrieves an object from memory would avoid such issues.

**Grasping module limitations:** One potential failure mode of our system is that our manipulation is performed by executing the outputs of a pre-trained model-generated grasps that are predicted based on a single RGB-D image, with a model that wasn’t designed for the Hello Robot: Stretch gripper.

As a result, sometimes such grasps are unreliable or unrealistic, as shown in Figure 8. There are cases where the proposed grasp is infeasible given the robot joint limits, or is simply too far from the robot body. Development of better heuristics will let us sample better grasps for a given object. In some other cases, the model generates a good grasp pose, but as the robot is executing the grasping primitive, it collides with some minor environment obstacle. Since we do not plan the grasp trajectory, and instead try to apply the same grasp trajectory in every case, some such failures are inevitable. Better grasping models that generates a grasp trajectory as well as a pose may solve such issues. Finally, our grasping module struggles with flat objects categorically, like chocolate bars and books, since it’s difficult to grasp them off a surface with a two-fingered gripper.

**Robot hardware limitations:** While our robot of choice, a



Hello Robot: Stretch, is able to pick-and-drop a number of objects, there are certain hardware limitations that determines what the robot can and cannot manipulate. For example, the robot has a 1 kg (2 lbs) payload limit when the arm is fully extended, and as such our method is unable to move objects like a full dish soap container. Similarly, objects that are far from navigable floor space, such as in the middle of a bed, or on high places, is difficult for the robot to reach because of the reach limits of the arm. Finally, in some situations, the robot hardware or the RealSense camera can become miscalibrated over time, especially during continuous testing in homes. This miscalibration can lead to error since the manipulation module requires hand-eye coordination in the robot.

## IV. RELATED WORKS

### A. Vision-Language models for robotic navigation

Early applications of pre-trained open-knowledge models has been in open-vocabulary navigation. Navigating to various objects is an important task which has been looked at in a wide range of previous works [25, 31, 39], as well as in the context of longer pick-and-place tasks [40, 41]. However, these methods have generally been applied to relatively small numbers of objects [42]. Recently, Objaverse [43] has shown navigation to thousands of object types, for example, but much of this work has been restricted to simulated or highly controlled environments.

The early work addressing this problem builds upon representations derived from pre-trained vision language models, such as SemAbs [44], CLIP-Fields [27], VLMaps [45], NLMap-SayCan [46], and later, ConceptFusion [47] and LERF [30]. Most of these models show object localization in pre-mapped scenes, while CLIP-Fields, VLMaps, and NLMap-SayCan show integration with real robots for indoor navigation tasks. USA-Nets [26] extends this task to include an affordance model, navigating with open-vocabulary queries while doing object avoidance. ViNT [48] proposes a foundation model for robotic navigation which can be applied to vision-language navigation problems. More recently, GOAT [31] was proposed as a modular system for “going to anything” and navigating to any object in any environment. ConceptGraphs [49] proposed an open scene representation capable of handling complex queries using LLMs and creating a scene graph.

### B. Pretrained robot manipulation models

While humans can frequently look at objects and immediately know how to grasp it, such grasping knowledge is not easily accessible to robots. Over the years, there has been many works that has focused on creating such a general robot grasp generation model [1, 50–55] for arbitrary objects and potentially cluttered scenes via learning methods. Our work focuses on more recent iterations of such methods [16, 19] that are trained on large grasping datasets [18, 38]. While these models only perform one task, namely grasping, they predict grasps across a large object surface and thus enable downstream complex, long-horizon manipulation tasks [20, 21, 56].

More recently, there is a set of general-purpose manipulation models moving beyond just grasping [57–61]. Some of these works perform general language-conditioned manipulation tasks, but are largely limited to a small set of scenes and objects. HACMan [62] demonstrates a larger range of object manipulation capabilities, focused on pushing and prodding. In the future, such models could expand the reach of our system.

### C. Open vocabulary robot systems

Many recent works have worked on language-enabled tasks for complex robot systems. Some examples include language conditioned policy learning [57, 63–65], learning goal-conditioned value functions [3, 66], and using large language models to generate code [67–69]. However, a fundamental difference remains between systems which aim to operate on arbitrary objects in an open-vocab manner, and systems where one can specify one among a limited number of goals or options using language. Consequently, Open-Vocabulary Mobile Manipulation has been proposed as a key challenge for robotic manipulation [25]. There has previously been efforts to build such a system [70, 71]. However, unlike such previous work, we try to build everything on an open platform and ensure our method can work without having to re-train anything for a novel home. Recently, UniTeam [23] won the 2023 HomeRobot OVMM Challenge [22] with a modular system doing pick-and-place to arbitrary objects, with a zero-shot generalization requirement similar to ours.

In parallel, recently, there have been a number of papers doing open-vocabulary manipulation using GPT or especially GPT4 [35]. GPT4V can be included in robot task planning frameworks and used to execute long-horizon robot tasks, including ones from human demonstrations [72]. ConceptGraphs [49] is a good recent example, showing complex object search, planning, and pick-and-place capabilities to open-vocabulary objects. SayPlan [73] also shows how these can be used together with a scene graph to handle very large, complex environments, and multi-step tasks; this work is complementary to ours, as it doesn’t handle how to implement pick and place.

## V. LIMITATIONS, OPEN PROBLEMS AND REQUEST FOR RESEARCH

While our method shows significant success in completely novel home environments, it also shows many places where such methods can improve. In this section, we discuss a few of such potential improvement in the future.

### A. Live semantic memory and obstacle maps

All the current semantic memory modules and obstacle map builders build a static representation of the world, without a good way of keeping it up-to-date as the world changes. However, homes are dynamic environments, with many small changes over the day every day. Future research that can build a dynamic semantic memory and obstacle map would unlock potential for continuous application of such pick-and-drop methods in a novel home out of the box.

### B. Grasp plans instead of proposals

Currently, the grasping module proposes generic grasps without taking the robot’s body and dynamics into account. Similarly, given a grasp pose, often the open loop grasping trajectory collides with environmental obstacles, which can be easily improved by using a module to generate grasp plans rather than grasp poses only.

### C. Improving interactivity between robot and user

One of the major causes of failure in our method is in navigation: where the semantic query is ambiguous and the intended object is not retrieved from the semantic memory. In such ambiguous cases, interaction with the user would go a long way to disambiguate the query and help the robot succeed more often.

### D. Detecting and recovering from failure

Currently, we observe a multiplicative error accumulation between our modules: if any of our independent components fail, the entire process fails. As a result, even if our modules each perform independently at or above 80% success rate, our final success rate can still be below 60%. However, with better error detection and retrying algorithms, we can recover from much more single-stage errors, and similarly improve our overall success rate [23].

### E. Robustifying robot hardware

While Hello Robot - Stretch [74] is an affordable and portable platform on which we can implement such an open-home system for arbitrary homes, we also acknowledge that with robust hardware such methods may have vastly enhanced capacity. Such robust hardware may enable us to reach high and low places, and pick up heavier objects. Finally, improved robot odometry will enable us to execute much more finer grasps than is possible today.

#### ACKNOWLEDGMENTS

NYU authors are supported by grants from Amazon, Honda, and ONR award numbers N00014-21-1-2404 and N00014-21-1-2758. NMS is supported by the Apple Scholar in AI/ML Fellowship. LP is supported by the Packard Fellowship. Our utmost gratitude goes to our friends and colleagues who helped us by hosting our experiments in their homes. Finally, we thank Jay Vakil, Siddhant Halder, Paula Pascual and Ulyana Piterbarg for valuable feedback and conversations.

#### REFERENCES

- [1] Lerrel Pinto and Abhinav Gupta. Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours. *ICRA*, 2016.
- [2] Sergey Levine, Peter Pastor, Alex Krizhevsky, Julian Ibarz, and Deirdre Quillen. Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. *The International journal of robotics research*, 37(4-5):421–436, 2018.
- [3] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn,

Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, et al. Do as I can, not as I say: Grounding language in robotic affordances. *Conference on Robot Learning (CoRL)*, 2022.

- [4] Nur Muhammad Mahi Shafiqullah, Anant Rai, Haritheja Etukuru, Yiqian Liu, Ishan Misra, Soumith Chintala, and Lerrel Pinto. On bringing robots home, 2023.
- [5] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- [6] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.
- [7] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In *European Conference on Computer Vision*, pages 350–368. Springer, 2022.
- [8] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, Xiao Wang, Xiaohua Zhai, Thomas Kipf, and Neil Houlsby. Simple open-vocabulary object detection with vision transformers. In *European Conference on Computer Vision*, pages 728–755. Springer, 2022.
- [9] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning (ICML)*, volume 139, pages 8748–8763, 2021.
- [10] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, pages 3195–3204, 2019.
- [11] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning, 2022.
- [12] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding DINO: Marrying DINO with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023.

- [13] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023.
- [14] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI Blog*, 2019.
- [15] Theophile Gervet, Soumith Chintala, Dhruv Batra, Jitendra Malik, and Devendra Singh Chaplot. Navigating to objects in the real world. *Science Robotics*, 8(79):eadf6991, 2023.
- [16] Martin Sundermeyer, Arsalan Mousavian, Rudolph Triebel, and Dieter Fox. Contact-graspnet: Efficient 6-dof grasp generation in cluttered scenes. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 13438–13444. IEEE, 2021.
- [17] Jeffrey Mahler, Jacky Liang, Sherdil Niyaz, Michael Laskey, Richard Doan, Xinyu Liu, Juan Aparicio Ojea, and Ken Goldberg. Dex-net 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics. *arXiv preprint arXiv:1703.09312*, 2017.
- [18] Hao-Shu Fang, Chenxi Wang, Minghao Gou, and Cewu Lu. Graspnet-1billion: a large-scale benchmark for general object grasping. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11444–11453, 2020.
- [19] Hao-Shu Fang, Chenxi Wang, Hongjie Fang, Minghao Gou, Jirong Liu, Hengxu Yan, Wenhai Liu, Yichen Xie, and Cewu Lu. Anygrasp: Robust and efficient grasp perception in spatial and temporal domains. *IEEE Transactions on Robotics*, 2023.
- [20] Ankit Goyal, Arsalan Mousavian, Chris Paxton, Yu-Wei Chao, Brian Okorn, Jia Deng, and Dieter Fox. Ifor: Iterative flow minimization for robotic object rearrangement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14787–14797, 2022.
- [21] Weiyu Liu, Tucker Hermans, Sonia Chernova, and Chris Paxton. Structdiffusion: Object-centric diffusion for semantic rearrangement of novel objects. *arXiv preprint arXiv:2211.04604*, 2022.
- [22] Sriram Yenamandra, Arun Ramachandran, Mukul Khanna, Karmesh Yadav, Devendra Singh Chaplot, Gunjan Chhablani, Alexander Clegg, Theophile Gervet, Vidhi Jain, Ruslan Partsey, Ram Ramrakhya, Andrew Szot, Tsung-Yen Yang, Aaron Edsinger, Charlie Kemp, Binit Shah, Zsolt Kira, Dhruv Batra, Roozbeh Mottaghi, Yonatan Bisk, and Chris Paxton. The homerobot open vocab mobile manipulation challenge. In *Thirty-seventh Conference on Neural Information Processing Systems: Competition Track*, 2023.
- [23] Andrew Melnik, Michael Büttner, Leon Harz, Lyon Brown, Gora Chand Nandi, Arjun PS, Gaurav Kumar Yadav, Rahul Kala, and Robert Haschke. Uniteam: Open vocabulary mobile manipulation challenge. *arXiv preprint arXiv:2312.08611*, 2023.
- [24] Luca Medeiros. Lang segment anything. <https://github.com/luca-medeiros/lang-segment-anything>, 2023.
- [25] Sriram Yenamandra, Arun Ramachandran, Karmesh Yadav, Austin Wang, Mukul Khanna, Theophile Gervet, Tsung-Yen Yang, Vidhi Jain, Alexander William Clegg, John Turner, et al. Homerobot: Open-vocabulary mobile manipulation. *arXiv preprint arXiv:2306.11565*, 2023.
- [26] Benjamin Bolte, Austin Wang, Jimmy Yang, Mustafa Mukadam, Mrinal Kalakrishnan, and Chris Paxton. Usanet: Unified semantic and affordance representations for robot memory. *arXiv preprint arXiv:2304.12164*, 2023.
- [27] Nur Muhammad Mahi Shafullah, Chris Paxton, Lerrel Pinto, Soumith Chintala, and Arthur Szlam. Clip-fields: Weakly supervised semantic fields for robotic memory. *arXiv preprint arXiv:2210.05663*, 2022.
- [28] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollar, and Ross Girshick. Segment anything. In *ICCV*, pages 4015–4026, October 2023.
- [29] David Rozenberszki, Or Litany, and Angela Dai. Language-grounded indoor 3d semantic segmentation in the wild, 2022.
- [30] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lurf: Language embedded radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19729–19739, 2023.
- [31] Matthew Chang, Theophile Gervet, Mukul Khanna, Sriram Yenamandra, Dhruv Shah, So Yeon Min, Kavit Shah, Chris Paxton, Saurabh Gupta, Dhruv Batra, et al. Goat: Go to any thing. *arXiv preprint arXiv:2311.06430*, 2023.
- [32] Chenguang Huang, Oier Mees, Andy Zeng, and Wolfram Burgard. Audio visual language maps for robot navigation. *arXiv preprint arXiv:2303.07522*, 2023.
- [33] Santiago Garrido, Luis Moreno, Mohamed Abderrahim, and Fernando Martin. Path planning for mobile robot navigation using voronoi diagram and fast marching. In *2006 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 2376–2381. IEEE, 2006.
- [34] Sudeep Dasari, Abhinav Gupta, and Vikash Kumar. Learning dexterous manipulation from exemplar object trajectories and pre-grasps, 2023.
- [35] OpenAI. GPT-4 technical report, 2023.
- [36] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks, 2019.
- [37] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *European Conference on Computer Vision (ECCV)*, 65(1):99–106, 2020.
- [38] Clemens Eppner, Arsalan Mousavian, and Dieter Fox. Acronym: A large-scale grasp dataset based on simulation. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6222–6227. IEEE, 2021.
- [39] Arsalan Mousavian, Alexander Toshev, Marek Fišer, Jana Košecká, Ayzaan Wahid, and James Davidson. Visual representations for semantic target driven navigation. In

- 2019 *International Conference on Robotics and Automation (ICRA)*, pages 8846–8852. IEEE, 5 2019.
- [40] Valts Blukis, Chris Paxton, Dieter Fox, Animesh Garg, and Yoav Artzi. A persistent spatial semantic representation for high-level natural language instruction execution. In *Conference on Robot Learning*, pages 706–717. PMLR, 2022.
- [41] So Yeon Min, Devendra Singh Chaplot, Pradeep Ravikumar, Yonatan Bisk, and Ruslan Salakhutdinov. Film: Following instructions in language with modular methods. *arXiv preprint arXiv:2110.07342*, 2021.
- [42] Matt Deitke, Dhruv Batra, Yonatan Bisk, Tommaso Campari, Angel X Chang, Devendra Singh Chaplot, Changan Chen, Claudia Pérez D’Arpino, Kiana Ehsani, Ali Farhadi, et al. Retrospectives on the embodied ai workshop. *arXiv preprint arXiv:2210.06849*, 2022.
- [43] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13142–13153, 2023.
- [44] Huy Ha and Shuran Song. Semantic abstraction: Open-world 3d scene understanding from 2d vision-language models, 2022.
- [45] Chenguang Huang, Oier Mees, Andy Zeng, and Wolfram Burgard. Visual language maps for robot navigation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 10608–10615. IEEE, 2023.
- [46] Boyuan Chen, Fei Xia, Brian Ichter, Kanishka Rao, Keerthana Gopalakrishnan, Michael S. Ryoo, Austin Stone, and Daniel Kappler. Open-vocabulary queryable scene representations for real world planning. In *arXiv preprint arXiv:2209.09874*, 2022.
- [47] Krishna Murthy Jatavallabhula, Alihusein Kuwajerwala, Qiao Gu, Mohd Omama, Tao Chen, Shuang Li, Ganesh Iyer, Soroush Saryazdi, Nikhil Keetha, Ayush Tewari, et al. Conceptfusion: Open-set multimodal 3d mapping. *arXiv preprint arXiv:2302.07241*, 2023.
- [48] Dhruv Shah, Ajay Sridhar, Nitish Dashora, Kyle Stachowicz, Kevin Black, Noriaki Hirose, and Sergey Levine. ViNT: A Foundation Model for Visual Navigation. In *7th Annual Conference on Robot Learning (CoRL)*, 2023.
- [49] Qiao Gu, Alihusein Kuwajerwala, Sacha Morin, Krishna Murthy Jatavallabhula, Bipasha Sen, Aditya Agarwal, Corban Rivera, William Paul, Kirsty Ellis, Rama Chellappa, et al. Conceptgraphs: Open-vocabulary 3d scene graphs for perception and planning. *arXiv preprint arXiv:2309.16650*, 2023.
- [50] Abhinav Gupta, Adithyavairavan Murali, Dhiraj Prakashchand Gandhi, and Lerrel Pinto. Robot learning in homes: Improving generalization and reducing dataset bias. *Advances in Neural Information Processing Systems*, 31:9094–9104, 2018.
- [51] Jeffrey Mahler, Jacky Liang, Sherdil Niyaz, Michael Laskey, Richard Doan, Xinyu Liu, Juan Aparicio Ojea, and Ken Goldberg. Dex-Net 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics. In *Robotics: Science and Systems (RSS)*, 2017.
- [52] Jeffrey Mahler, Matthew Matl, Xinyu Liu, Albert Li, David Gealy, and Ken Goldberg. Dex-net 3.0: Computing robust robot vacuum suction grasp targets in point clouds using a new analytic model and deep learning, 2018.
- [53] Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. QT-Opt: Scalable deep reinforcement learning for vision-based robotic manipulation. *arXiv preprint arXiv:1806.10293*, 2018.
- [54] Yuzhe Qin, Rui Chen, Hao Zhu, Meng Song, Jing Xu, and Hao Su. S4g: Amodal single-view single-shot se(3) grasp detection in cluttered scenes, 2019.
- [55] Arsalan Mousavian, Clemens Eppner, and Dieter Fox. 6-dof grasping: Variational grasp generation for object manipulation, 2019.
- [56] I. Singh, V. Blukis, A. Mousavian, A. Goyal, D. Xu, J. Tremblay, and D. Fox. Progprompt: Generating situated robot task plans using large language models. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, page 11523, 2023.
- [57] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Perceiver-Actor: A multi-task transformer for robotic manipulation. In *CoRL*, pages 785–799. PMLR, 2023.
- [58] Priyam Parashar, Jay Vakil, Sam Powers, and Chris Paxton. Spatial-language attention policies for efficient robot learning. *arXiv preprint arXiv:2304.11235*, 2023.
- [59] Nur Muhammad Shafiullah, Zichen Cui, Ariuntuya Arty Altanzaya, and Lerrel Pinto. Behavior transformers: Cloning  $k$  modes with one stone. *Advances in neural information processing systems*, 35:22955–22968, 2022.
- [60] Zichen Jeff Cui, Yibin Wang, Nur Muhammad Mahi Shafiullah, and Lerrel Pinto. From play to policy: Conditional behavior generation from uncurated robot data, 2022.
- [61] Theophile Gervet, Zhou Xian, Nikolaos Gkanatsios, and Katerina Fragkiadaki. Act3d: 3d feature field transformers for multi-task robotic manipulation. In *Conference on Robot Learning*, pages 3949–3965. PMLR, 2023.
- [62] Wenxuan Zhou, Bowen Jiang, Fan Yang, Chris Paxton, and David Held. Learning hybrid actor-critic maps for 6d non-prehensile manipulation. *arXiv preprint arXiv:2305.03942*, 2023.
- [63] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. CLIPort: What and where pathways for robotic manipulation. In *CoRL*, pages 894–906. PMLR, 2022.
- [64] Corey Lynch, Mohi Khansari, Ted Xiao, Vikash Kumar, Jonathan Tompson, Sergey Levine, and Pierre Sermanet. Learning latent plans from play. In *CoRL*, pages 1113–1132. PMLR, 2020.
- [65] Corey Lynch and Pierre Sermanet. Language conditioned imitation learning over unstructured data. *Robotics:*

*Science and Systems*, 2021.

- [66] Wenlong Huang, Chen Wang, Ruohan Zhang, Yunzhu Li, Jiajun Wu, and Li Fei-Fei. VoxPoser: Composable 3D value maps for robotic manipulation with language models. In *CoRL*, 2023.
- [67] Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as Policies: Language model programs for embodied control. In *icra*, pages 9493–9500. IEEE, 2023.
- [68] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv: Arxiv-2305.16291*, 2023.
- [69] Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. ProgPrompt: Generating situated robot task plans using large language models. In *ICRA*, pages 11523–11530. IEEE, 2023.
- [70] Naoki Yokoyama, Alex Clegg, Joanne Truong, Eric Undersander, Tsung-Yen Yang, Sergio Arnaud, Sehoon Ha, Dhruv Batra, and Akshara Rai. ASC: Adaptive skill coordination for robotic mobile manipulation. *arXiv preprint arXiv:2304.00410*, 2023.
- [71] Austin Stone, Ted Xiao, Yao Lu, Keerthana Gopalakrishnan, Kuang-Huei Lee, Quan Vuong, Paul Wohlhart, Sean Kirmani, Brianna Zitkovich, Fei Xia, Chelsea Finn, and Karol Hausman. Open-world object manipulation using pre-trained vision-language model. In *arXiv preprint*, 2023.
- [72] Naoki Wake, Atsushi Kanehira, Kazuhiro Sasabuchi, Jun Takamatsu, and Katsushi Ikeuchi. Gpt-4v(ision) for robotics: Multimodal task planning from human demonstration. *arXiv preprint arXiv:2311.12015*, 2023.
- [73] Krishan Rana, Jesse Haviland, Sourav Garg, Jad Abou-Chakra, Ian Reid, and Niko Suenderhauf. Sayplan: Grounding large language models using 3d scene graphs for scalable task planning. *arXiv preprint arXiv:2307.06135*, 2023.
- [74] Charles C Kemp, Aaron Edsinger, Henry M Clever, and Blaine Matulevich. The design of stretch: A compact, lightweight mobile manipulator for indoor human environments. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 3150–3157. IEEE, 2022.

## APPENDIX A SCANNET200 TEXT QUERIES

To detect objects in a given home environment using OWL-ViT, we use the Scannet200 labels. The full label set is here: ['shower head', 'spray', 'inhaler', 'guitar case', 'plunger', 'range hood', 'toilet paper dispenser', 'adapter', 'soy sauce', 'pipe', 'bottle', 'door', 'scale', 'paper towel', 'paper towel roll', 'stove', 'mailbox', 'scissors', 'tape', 'bathroom stall', 'chopsticks', 'case of water bottles', 'hand sanitizer', 'laptop', 'alcohol disinfection', 'keyboard', 'coffee maker', 'light', 'toaster', 'stuffed animal', 'divider', 'clothes dryer', 'toilet seat cover dispenser', 'file cabinet', 'curtain', 'ironing board', 'fire extinguisher', 'fruit', 'object', 'blinds', 'container', 'bag', 'oven', 'body wash', 'bucket', 'cd case', 'tv', 'tray', 'bowl', 'cabinet', 'speaker', 'crate', 'projector', 'book', 'school bag', 'laundry detergent', 'mattress', 'bathtub', 'clothes', 'candle', 'basket', 'glass', 'face wash', 'notebook', 'purse', 'shower', 'power outlet', 'trash bin', 'paper bag', 'water dispenser', 'package', 'bulletin board', 'printer', 'windowsill', 'disinfecting wipes', 'bookshelf', 'recycling bin', 'headphones', 'dresser', 'mouse', 'shower gel', 'dustpan', 'cup', 'storage organizer', 'vacuum cleaner', 'fireplace', 'dish rack', 'coffee kettle', 'fire alarm', 'plants', 'rag', 'can', 'piano', 'bathroom cabinet', 'shelf', 'cushion', 'monitor', 'fan', 'tube', 'box', 'blackboard', 'ball', 'bicycle', 'guitar', 'trash can', 'hand sanitizers', 'paper towel dispenser', 'whiteboard', 'bin', 'potted plant', 'tennis', 'soap dish', 'structure', 'calendar', 'dumbbell', 'fish oil', 'paper cutter', 'ottoman', 'stool', 'hand wash', 'lamp', 'toaster oven', 'music stand', 'water bottle', 'clock', 'charger', 'picture', 'basketball', 'sink', 'microwave', 'screwdriver', 'kitchen counter', 'rack', 'apple', 'washing machine', 'suitcase', 'ladder', 'ping pong ball', 'window', 'dishwasher', 'storage container', 'toilet paper holder', 'coat rack', 'soap dispenser', 'refrigerator', 'banana', 'counter', 'toilet paper', 'mug', 'marker pen', 'hat', 'aerosol', 'luggage',

'poster', 'bed', 'cart', 'light switch', 'backpack', 'power strip', 'baseball', 'mustard', 'bathroom vanity', 'water pitcher', 'closet', 'couch', 'beverage', 'toy', 'salt', 'plant', 'pillow', 'broom', 'pepper', 'muffins', 'multivitamin', 'towel', 'storage bin', 'nightstand', 'radiator', 'telephone', 'pillar', 'tissue box', 'vent', 'hair dryer', 'ledge', 'mirror', 'sign', 'plate', 'tripod', 'chair', 'kitchen cabinet', 'column', 'water cooler', 'plastic bag', 'umbrella', 'doorframe', 'paper', 'laundry hamper', 'food', 'jacket', 'closet door', 'computer tower', 'stairs', 'keyboard piano', 'person', 'table', 'machine', 'projector screen', 'shoe'].

## APPENDIX B SAMPLE OBJECTS FROM OUR TRIALS

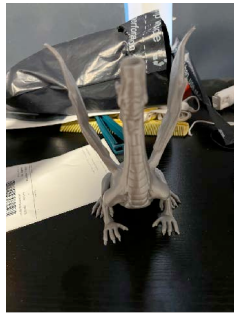
During our experiments, we tried to sample objects that can plausibly be manipulated by the Hello Robot: Stretch gripper from the home environments. As a result, OK-Robot encountered a large variety of objects with different shapes and visual features. A subsample of such objects are presented in the Figures 9, 10.

## APPENDIX C LIST OF HOME EXPERIMENTS

A full list of experiments in homes can be found in Table I.



Arm smartphone holder



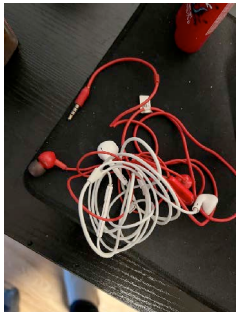
Gray toy dragon



Toy plant



White shirt



Tangled earphones



Playing cards



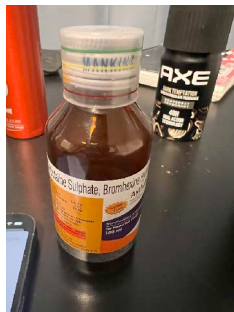
Blue gloves



Toy cactus



Toy grapes



Medicine bottles



Grey rag



Blue hair oil bottle



Blue pretzel pack



Toothpaste



White pretzel



Blue body wash

Fig. 9: Sample objects on our home experiments, sampled from each home environment, which OK-Robot was able to pick and drop successfully.



Purple strap



Yellow ginger paste packet



Blue bag



Steel wool



Translucent grey cup



Black face wash



Gold wrapped chocolate



Black head band



Blue eyeglass case



Fluffy headbands



Yogurt drinks



Lotion pump



Blue hair gel tube



Brown trail mix bag



White Apple bag



Small hand sanitizer

Fig. 10: Sample objects on our home experiments, sampled from each home environment, which OK-Robot **failed** to pick up successfully.



TABLE I: A list of all tasks in the home environments, along with their categories and success rates out of 10 trials.

Pick object	Place location	Result
Home 1		
Cleanup level: none		
silver cup	white table	Success
blue eye glass case	chair	Success
printed paper cup, coffee cup [white table]	—	Manipulation failure
small red and white medication	Chair	Success
power adapter	Grey Bed	Success
wrapped paper	—	Navigation failure
blue body wash	study table	Success
blue air spray	white table	Success
black face wash	—	Manipulation failure
yellow face wash	chair	Success
body spray	—	Navigation failure
small hand sanitizer	—	Manipulation failure
blue inhaler device(window)	white table	Success
inhaler box(window)	dust bin	Success
multivitamin container	—	Navigation failure
red towel	white cloth bin (air conditioner)	Success
white shirt	white cloth bin (air conditioner)	Success
Cleanup level: low		
silver cup	white table	Success
blue eye glass case	—	Navigation failure
printed paper cup, coffee cup [white table]	dust bin	Success
small red and white medication	Chair	Success
power adapter	—	Navigation failure
blue body wash	white table	Success
blue air spray	white table	Success
yellow face wash	white table	Success
small hand sanitizer	study table	Success
blue inhaler device(window)	—	Manipulation failure
inhaler box(window)	dust bin	Success
red towel	white cloth bin(air conditioner)	Success
white shirt	white cloth bin(air conditioner)	Success
Cleanup level: high		
silver cup	white table	Success
printed paper cup, coffee cup [white table]	dust bin	Success
blue body wash	white table	Success
blue air spray	white table	Success
yellow face wash	—	Manipulation failure
small hand sanitizer	—	Manipulation failure
inhaler box(window)	dust bin	Success
white shirt	white cloth bin(air conditioner)	Success
Home 2		
Cleanup level: None		
fanta can	dust bin	Success
tennis ball	small red shopping bag	Success
black head band [bed]	—	Manipulation failure
purple shampoo bottle	white rack	Success
toothpaste	small red shopping bag	Success
Continued on the next page		

Pick object	Place location	Result
orange packaging	dust bin	Success
green hair cream jar [white rack]	—	Navigation failure
green detergent pack [white rack]	white table	Success
blue moisturizer [white rack]	—	Navigation failure
green plastic cover	—	Navigation failure
storage container	—	Manipulation failure
blue hair oil bottle	white rack	Success
blue pretzels pack	white rack	Success
blue hair gel tube	—	Manipulation failure
red bottle [white rack]	brown desk	Success
blue bottle [air conditioner]	white cloth bin(air conditioner)	Success
wallet	—	Manipulation failure

#### Cleanup level: low

fanta can	black trash can	Success
tennis ball	red target bag	Success
black head band [bed]	red target bag	Success
purple shampoo bottle	red target bag	Success
toothpaste	red target bag	Success
orange packaging	black trash can	Success
green detergent pack [white rack]	—	Manipulation failure
blue moisturizer [white rack]	—	Navigation failure
blue hair oil bottle	white rack	Success
blue pretzels pack	white rack	Success
wallet	—	Manipulation failure

#### Cleanup level: high

fanta can	black trash can	Success
purple shampoo bottle	small red shopping bag	Success
orange packaging	black trash can	Success
blue moisturizer [white rack]	white rack	Success
blue hair oil bottle	—	Manipulation failure
blue hair gel tube	dust bin	Success
red bottle [white rack]	target bag	Placing failure
blue bottle [air conditioner]	white cloth bin(air conditioner)	Success

#### Home 3

#### Cleanup level: none

apple	white plate	Success
ice cream	white and green bag	Success
green lime juice bottle	red basket	Success
yellow packet	—	Manipulation failure
red packet	—	Manipulation failure
orange can	card board box	Success
cooking oil bottle	—	Manipulation failure
pasta sauce	—	Manipulation failure
orange box [stove]	—	Manipulation failure
green bowl	sink	Success
washing gloves	green bag [card board box]	Success
small oregano bottle	red basket	Success
yellow noodles packet [stove]	red basket	Success
blue dish wash bottle	card board box	Success
scrubber	—	Navigation failure
dressing salad bottle	—	Navigation failure

Continued on the next page

Pick object	Place location	Result
Cleanup level: low		
apple	white plate	Success
ice cream	red basket	Success
green lime juice bottle	red basket	Success
yellow packet	green bag	Success
red packet	_____	Manipulation failure
orange can	card board box	Success
cooking oil bottle	marble surface [red basket]	Success
green bowl	_____	Manipulation failure
washing gloves	sink	Success
small oregano bottle	red basket	Success
yellow noodles packet [stove]	_____	Manipulation failure
blue dish wash bottle	card board box	Success
Cleanup level: high		
apple	white plate	Success
ice cream	red basket	Success
green lime juice bottle	red basket	Success
orange can	card board box	Success
cooking oil bottle	_____	Manipulation failure
washing gloves	sink	Success
small oregano bottle	red basket	Success
yellow noodles packet [stove]	red basket	Success
blue dish wash bottle	card board box	Success
Home 4		
Cleanup level: none		
pepsi	black chair	Success
birdie	cloth bin	Success
black hat	_____	Navigation failure
owl like wood carving	bed	Success
red inhaler	_____	Manipulation failure
black sesame seeds	_____	Manipulation failure
banana	_____	Manipulation failure
loose-leaf herbal tea jar	black chair	Success
red pencil sharpener	_____	Navigation failure
fast-food French fries container	blue shopping bag [metal drying rack]	Placing failure
milk	plastic storage drawer unit	Success
socks[bed]	_____	Navigation failure
purple gloves	_____	Manipulation failure
target bag	cloth bin	Success
muffin	grey bed	Success
tissue paper	table	Success
grey eyeglass box	_____	Manipulation failure
Cleanup level: low		
pepsi	basket	Success
birdie	white drawer	Success
owl like wood carving	_____	Navigation failure
red inhaler	plastic storage drawer unit	Success
black sesame seeds	bed	Success
loose-leaf herbal tea jar	table	Success
fast-food French fries container	chair	Success
Continued on the next page		

Pick object	Place location	Result
milk	chair	Success
purple gloves	basket	Success
target bag	basket	Placing failure
muffin	table	Success
tissue paper	—	Manipulation failure
grey eyeglass box	—	Navigation failure

Cleanup level: high

pepsi	basket	Success
birdie	bed	Success
red inhaler	plastic storage drawer unit	Success
black sesame seeds	desk	Success
banana	—	Manipulation failure
loose-leaf herbal tea jar	—	Manipulation failure
milk	chair	Success
purple gloves	basket	Success
target bag	basket	Success
muffin	bed	Success

Home 5

Cleanup level: none

tiger balm topical ointment	—	Navigation failure
pink shampoo	trader joes shapping bag	Success
aveeno sunscreen protective lotion	trader joes shapping bag	Success
small yellow nozzle spray	—	Manipulation failure
black hair care spray	—	Manipulation failure
green hand sanitizer	—	Manipulation failure
white hand sanitizer	—	Navigation failure
white bowl [ketchup]	black sofa chair	Success
blue bowl	—	Manipulation failure
blue sponge	trader joes shapping bag	Success
ketchup	—	Manipulation failure
white salt	—	Manipulation failure
black pepper	black drawer	Success
blue bottle	—	Navigation failure
purple light bulb box	trader joes shopping bag	Success
white plastic bag	bed	Success
rag	white rack	Success

Cleanup level: low

pink shampoo	—	Navigation failure
aveeno sunscreen protective lotion	—	Manipulation failure
small yellow nozzle spray	—	Manipulation failure
white bowl [ketchup]	black sofa chair	Success
blue sponge	bed	Success
ketchup	trader joes shopping bag	Success
white salt	trader joes shopping bag	Success
black pepper	—	Navigation failure
blue bottle	black sofa chair	Success
purple light bulb box	—	Manipulation failure
rag	white rack	Success

Cleanup level: high

pink shampoo	trader joes shopping bag	Success
--------------	--------------------------	---------

Continued on the next page

Pick object	Place location	Result
green hand sanitizer	black sofa chair	Success
white bowl [ketchup]	_____	Manipulation failure
blue sponge	bed	Success
ketchup	black drawer	Success
white salt	white drawer	Success
purple light bulb box	trader joes shopping bag	Success
rag	black sofa chair	Success

**Home 6**  
Cleanup level: none

translucent grey cup	_____	Manipulation failure
green mouth spray box	stove	Success
green eyeglass container	chair	Success
blue bag	_____	Manipulation failure
black burn ointment box	_____	Navigation failure
white vitamin bottle	_____	Navigation failure
McDonald's paper bag	stove	Success
purple medicine packaging	chair	Success
grey rag	sink	Success
sparkling water can [sink]	countertop	Success
gold wrapped chocolate	_____	Manipulation failure
lemon tea carton	table	Success
metallic golden beverage can	table	Success
red bottle	table	Success
tea milk bottle	_____	Navigation failure
nyu water bottle [sink]	table	Success
white hand wash	_____	Navigation failure

Cleanup level: low

translucent grey cup	_____	Navigation failure
green mouth spray box	_____	Manipulation failure
blue bag	brown box	Success
black burn ointment box	brown box	Success
McDonald's paper bag	_____	Navigation failure
grey rag	sink	Success
sparkling water can [sink]	chair	Success
lemon tea carton	stove	Success
metallic golden beverage can	_____	Navigation failure
red bottle	brown box	Success
nyu water bottle [sink]	table	Success
white hand wash	sink	Success

Cleanup level: high

blue bag	brown box	Success
black burn ointment box	_____	Manipulation failure
grey rag	sink	Success
sparkling water can [sink]	chair	Success
lemon tea carton	table	Success
metallic golden beverage can	stove	Success
red bottle	_____	Navigation failure
nyu water bottle [sink]	_____	Manipulation failure
white hand wash	_____	Manipulation failure

**Home 7**  
Continued on the next page

Pick object	Place location	Result
Cleanup level: none		
blue plastic bag roll	_____	Navigation failure
green bag	basket[window]	Success
toy cactus	desk	Success
toy van	chair	Success
brown medical bandage	chair	Success
power adapter	_____	Navigation failure
red herbal tea	brown cardboard box	Success
apple juice box	brown cardboard box	Success
paper towel	blue cardboard box	Success
toy bear	bed blanket	Success
yellow ball	bed blanket	Success
black pants	basket[window]	Success
purple water bottle	desk	Success
blue eyeglass case	_____	Manipulation failure
brown toy monkey	_____	Navigation failure
blue hardware box [table]	blue cardboard box	Success
green zandu balm container	blue cardboard box	Success

Cleanup level: low		
green bag	basket	Success
toy cactus	basket	Success
toy van	chair	Success
brown medical bandage	_____	Manipulation failure
red herbal tea	brown box	Success
apple juice box	brown box	Success
paper towel	basket	Success
toy bear	desk	Success
purple water bottle	desk	Success
blue eyeglass case	_____	Manipulation failure
green zandu balm container	blue cardboard box	Success

Cleanup level: high		
green bag	stool [window]	Success
toy cactus	table	Success
toy van	white basket	Success
red herbal tea	brown cardboard box	Success
apple juice box	brown cardboard box	Success
paper towel	blue cardboard box	Success
toy bear	white basket	Success
yellow ball	bed	Success
purple water bottle	black tote bag	Success
green zandu balm container	blue cardboard box	Success

Home 8		
Cleanup level: none		
cyan air spray	brown shelf [sink]	Success
blue gloves	kitchen sink	Success
blue peanut butter	black stove	Success
nutella	table	Success
green bag	brown shelf [sink]	Success
green bandage box	trash can	Success
green detergent	kitchen sink	Success
black 'red pepper sauce'	_____	Manipulation failure

Continued on the next page

Pick object	Place location	Result
red bag	chair	Success
black bag	chair	Success
red spray [brown shelf]	kitchen countertop	Success
steel wool	_____	Manipulation failure
white aerosol	trash can	Success
white pretzel	black stove	Success
purple crisp	kitchen countertop	Success
plastic bowl	_____	Manipulation failure
playing card	microwave	Success

Cleanup level: low

cyan air apray	chair	Success
blue gloves	sink	Success
blue peanut butter	_____	Navigation failure
green bag	brown shelf	Success
green bandage box	brown shopping bag	Success
green detergent	microwave	Success
red bag	_____	Manipulation failure
black bag	chair	Success
white aerosol	trash can	Success
white pretzel	black stove	Success
purple crisp	kitchen countertop	Success
plastic bowl	_____	Manipulation failure
playing card	microwave	Success

Cleanup level: high

cyan air apray	brown shelf [sink]	Success
blue gloves	stove	Success
blue peanut butter	black stove	Success
green bag	brown shelf [sink]	Success
green bandage box	microwave	Success
green detergent	_____	Manipulation failure
black bag	chair	Success
white aerosol	table	Success
purple crisp	chair	Success
playing card	microwave	Success

Home 9

Cleanup level: none

toy grapes	black laundry bag	Success
purple strap	_____	Manipulation failure
red foggy body spray	_____	Manipulation failure
arm smartphone holder	bed	Success
medicine bottle	_____	Manipulation failure
yogurt beverage	_____	Navigation failure
blue shaving cream can	_____	Navigation failure
blue cup	table	Success
purple tape	_____	Manipulation failure
black shoe brush	_____	Navigation failure
fluffy headband	_____	Manipulation failure
black water bottle	brown shopping bag	Placing failure
yellow eyeglass case	black chair	Success
paper cup	_____	Manipulation failure
lotion pump	_____	Manipulation failure

Continued on the next page

Pick object	Place location	Result
nasal spray	_____	Manipulation failure
plastic bag	trash basket	Success

Cleanup level: low

toy grapes	_____	Manipulation failure
red foggy body spray	brown paper bag	Success
arm smartphone holder	brown paper bag	Success
yogurt beverage	desk	Success
blue shaving cream can	black bag	Success
blue cup	black chair	Success
black shoe brush	_____	Manipulation failure
fluffy headband	_____	Navigation failure
black water bottle	folded chair	Success
nasal spray	_____	Navigation failure
plastic bag	trash basket	Success

Cleanup level: high

red foggy body spray	brown paper bag	Success
arm smartphone holder	_____	Manipulation failure
yogurt beverage	desk	Success
blue shaving cream can	black bag	Success
blue cup	black chair	Success
black water bottle	white bed	Success
nasal spray	folded chair	Success
plastic bag	trash basket	Success

Home 10

Cleanup level: none

grey toy dragon	bed	Success
purple body spray	_____	Manipulation failure
hand sanitizer	shelf	Success
toy plant	bed [shelf]	Success
brown trail mix bag	_____	Manipulation failure
hanging blue shirt	cloth bin	Success
white apple bag	_____	Manipulation failure
white and pink powder bottle	table	Success
cough syrup bottle	shelf	Success
tangled ear phones	office chair	Success
red deodrant stick[table]	chair	Success
black body spray	chair	Success
hair treatment medicine bottle	_____	Manipulation failure
green tea package	chair	Success
portable speaker [green tea package]	office chair	Success
wooden workout gripper	_____	Navigation failure
brown box	_____	Navigation failure
blue bulb adapter	office chair	Success
game controller	office chair	Success

Cleanup level: low

grey toy dragon	orange bag	Success
purple body spray	table	Success
hand sanitizer	_____	Navigation failure
toy plant	bed	Success
brown trail mix bag	_____	Manipulation failure

Continued on the next page



Pick object	Place location	Result
white and pink powder bottle	black chair [bed]	Success
cough syrup bottle	shelf [bed]	Success
red deodrant stick[table]	bed [rack]	Success
black body spray	rack [bed]	Placing failure
green tea package	orange bag	Success
brown box	black chair [bed]	Success
blue bulb adapter	_____	Manipulation failure
Cleanup level: high		
purple body spray	orange bag	Success
toy plant	bed	Success
white and pink powder bottle	_____	Navigation failure
cough syrup bottle	shelf [bed]	Success
red deodrant stick[table]	_____	Navigation failure
black body spray	black chair	Success
green tea package	table	Success
blue bulb adapter	shelf	Success

---



---