

ARNOLD: A Benchmark for Language-Grounded Task Learning With Continuous States in Realistic 3D Scenes

Ran Gong^{1*}, Jianguo Huang^{2,5*}, Yizhou Zhao¹, Haoran Geng^{2,5}, Xiaofeng Gao¹, Qingyang Wu⁴,
Wensi Ai¹, Ziheng Zhou¹, Demetri Terzopoulos¹, Song-Chun Zhu^{2,3,5}, Baoxiong Jia⁵, Siyuan Huang⁵

¹University of California, Los Angeles, ²Peking University, ³Tsinghua University,
⁴Columbia University, ⁵National Key Laboratory of General Artificial Intelligence, BIGAI

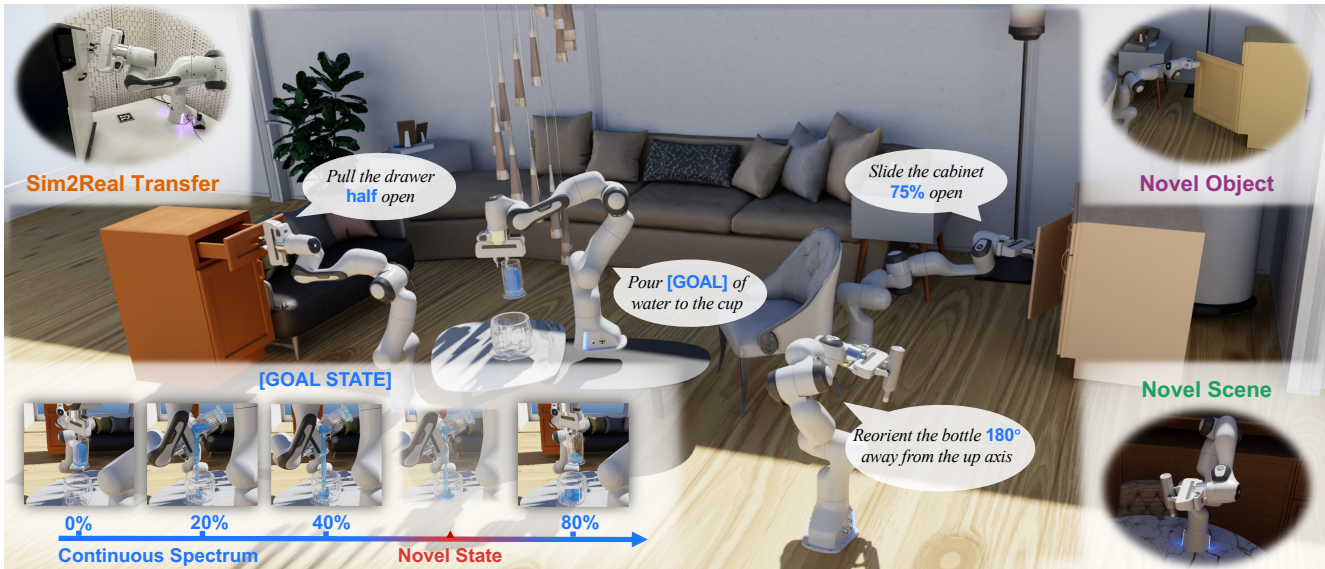


Figure 1. **The ARNOLD benchmark** for language-grounded task learning with **continuous states** in realistic 3D scenes. ARNOLD provides 8 tasks with their demonstrations for learning and a testbed for the generalization abilities of agents over (1) **novel goal states**, (2) **novel objects**, and (3) **novel scenes**.

Abstract

Understanding the continuous states of objects is essential for task learning and planning in the real world. However, most existing task learning benchmarks assume discrete (e.g., binary) object goal states, which poses challenges for the learning of complex tasks and transferring learned policy from simulated environments to the real world. Furthermore, state discretization limits a robot’s ability to follow human instructions based on the grounding of actions and states. To tackle these challenges, we present ARNOLD, a benchmark that evaluates language-grounded task learning with continuous states in realistic 3D scenes. ARNOLD is comprised of 8 language-conditioned tasks that involve understanding object states and learning policies for continuous goals. To promote language-instructed learning, we provide expert

demonstrations with template-generated language descriptions. We assess task performance by utilizing the latest language-conditioned policy learning models. Our results indicate that current models for language-conditioned manipulations continue to experience significant challenges in novel goal-state generalizations, scene generalizations, and object generalizations. These findings highlight the need to develop new algorithms that address this gap and underscore the potential for further research in this area.

Project website: <https://arnold-benchmark.github.io>.

1. Introduction

The ability to ground language is a crucial skill that has evolved over the course of human history, allowing people to learn and describe concepts, perform tasks, and communicate with one another. While recent developments have

* Equal contribution. Corresponding authors: Ran Gong, Jianguo Huang, Baoxiong Jia, and Siyuan Huang.

enabled the grounding of concepts in images [67, 40, 69, 21], the interaction with physical environments [10, 2, 84, 44, 62, 36, 35, 28, 26, 24, 87, 79, 9, 15], and language understanding in physical environments [50, 1, 65, 89, 77, 22, 80, 64], few researchers have investigated the grounding of actions in daily tasks [71, 90, 72, 81]. Given that humans can comprehend object status and relate language instructions to the physical environment, a pertinent question arises: *How can we imbue robotic systems with the same capacity to understand and execute language instructions in the physical world?*

Learning action grounding in daily activities is a challenging task that presents several non-trivial difficulties. Firstly, robotic tasks rely heavily on detailed scene understanding for successful execution. This includes the understanding of geometry information, layouts, and visual appearances. The various combinations of scene configurations, including novel appearances, objects, and spatial positions further exacerbate this challenge. Therefore, it is crucial for robotic systems to acquire generalizable skills that can be transferred to different domains and settings.

Secondly, humans possess an exquisite ability to understand desired goal states precisely. This ability allows us to effortlessly map from simple descriptions (*e.g.*, a cup half filled, a door fully opened, *etc.*) to the precise status of physical properties (*e.g.*, half the volume, pulled to 180°, *etc.*). However, it is exceedingly challenging for robots to learn the precise goal state from abstracted language instructions, especially when referring to an implicit range of continuous object states (*e.g.*, a bit of coffee, slightly open, *etc.*) [42, 30, 57]. As a result, there is an urgent need for robot systems to maintain a mapping from language instructions to precise goal states in a continuous world.

A necessary first step toward tackling these challenges is to develop realistic robot simulation systems that enable language-grounded learning. Indeed, notable recent advances in simulated environments have facilitated grounded task learning [8, 73, 58, 90, 54]. Despite the impressive progress, these benchmarks suffer from several limitations that hinder the effective operation of robots in the real world: (1) They typically assume that tasks are performed in simple and clean environments, rather than in scenes that are spatially occupied by clutter and visually disturbed by diverse textured backgrounds [43, 46, 90, 72]. (2) They assume discrete (*e.g.*, binary) object states and perfect motor control that ignore the low-level geometry and dynamics of objects [76, 74, 16] and, consequently, they do not attempt in-depth physical state understanding or fine-grained manipulation skills. (3) These benchmarks do not ground instructions to precise states [90, 71], thus neglecting the challenging problem of grounding language to object states.

To address these critical challenges of language-grounded robot task learning, we introduce a new benchmark,

ARNOLD, for grounding task instructions to precise robot actions and object states in realistic natural scenes (Fig. 1). Specifically, we leverage a highly accurate physics simulation engine to create eight challenging robot manipulation tasks that include continuous robot motion, friction-based grasping, and a variety of object state manipulations. Each task is associated with a set of goals sampled from a continuous range of object states and their corresponding detailed task descriptions in human language form. We further provide plentiful demonstrations of each task with trajectories generated with a template-based planner for robot learning.

To provide an in-depth evaluation of language-grounded task learning, we complement prior research with an evaluation that targets the ability of agents to generalize learned language-grounded skills to unseen scenarios, including novel scenes, novel objects, and our featured novel goal states. We have meticulously curated a collection of 40 distinctive objects and 20 scenes from open-source datasets [85, 41, 18] and designed data splits for evaluating different aspects of agents' generalization ability in language-grounded task learning. Furthermore, we provide thorough experimental analyses and show that state-of-the-art language-conditioned manipulation models still suffer with regard both to grounding and generalization. Additionally, we show that state modeling is crucial for tasks in ARNOLD through carefully designed ablation studies.

In summary, ARNOLD makes the following contributions:

- A **realistic 3D interactive environment** with diverse scenes, objects, and **continuous object states**, facilitating the learning and evaluation of precise robot manipulation.
- A systematic benchmark comprising eight challenging **language-grounded robotic tasks** and evaluation splits for different aspects of **skill generalization**.
- **Extensive experiments and analyses** of state-of-the-art language-conditioned manipulation models, revealing their strengths and weaknesses in promoting future research on language-grounded task learning.

2. Related Work

Simulators for Embodied AI. Significant progress has recently been made in developing simulators for training and evaluating AI agents to perform indoor household activities [44, 17, 19, 6, 11, 5]. To mitigate complexity, most of these simulators make simplifications about world states and actions, abstracting robot manipulation into symbolic planning in discrete action [66, 41] and state spaces. However, agents trained in such settings are unaware of the relationship between actions and the geometries and dynamics of objects, therefore limiting their abilities in real-world scenarios. Recent efforts have gradually transitioned to continuous action spaces, but they still make some simplifications. For example, grasping is often simplified by attaching a nearby object to the gripper [16, 74], or through contact [45, 37, 76].

Table 1. **Comparison with existing benchmarks.** ARNOLD features language-grounded robot control over continuous object states with a large number of demonstrations in photo-realistic scenes. ARNOLD also **leverages advanced physics simulations powered by PhysX 5.0 to simulate articulated bodies and fluids.** **Language:** Task goals are specified by human language instruction. **Multi-Camera:** Robot is equipped with multiple cameras. **Fluid:** Advanced fluid simulation. **Physics:** Realistic physics simulation with realistic grasping. ¹: RLbench-based benchmarks use simplified grasping. **Continuous:** Object state and goal state are continuous. **Scene:** Tasks are performed with a realistic scene background. **Robot:** Perform actions with real robots for all tasks. **R:** Rasterization. **RT:** RayTracing. **Flexible Material:** Easy to change materials and textures. **Generalization:** Systematic generalization test at different levels.

| Benchmark | Language | Multi Camera | Fluid | Physics | Continuous | Scenes | Robot | Rendering | Flexible Material | Generalization |
|------------------------|----------|--------------|-------|----------------|------------|--------|-------|-----------|-------------------|----------------|
| Alfred [73] | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | R | ✗ | ✗ |
| Maniskill [60, 29] | ✗ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | R | ✗ | ✓ |
| Calvin [58] | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ | ✓ | R | ✗ | ✓ |
| Behavior [45, 74] | ✗ | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ | RT | ✓ | ✗ |
| KitchenShift [86] | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | R | ✗ | ✓ |
| RLBench [37] | ✓ | ✓ | ✗ | ✓ ¹ | ✗ | ✗ | ✓ | R | ✗ | ✗ |
| Softgym [46] | ✗ | ✗ | ✓ | ✓ | ✓ | ✗ | ✗ | R | ✗ | ✗ |
| Orbit [59] | ✗ | ✗ | ✓ | ✓ | ✗ | ✓ | ✓ | RT | ✓ | ✗ |
| Vlmbench [90] | ✓ | ✓ | ✗ | ✓ ¹ | ✓ | ✗ | ✓ | R | ✗ | ✗ |
| Ravens [88, 71] | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ | R | ✗ | ✓ |
| Habitat HAB [76] | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | R | ✗ | ✓ |
| TDW Transport [19, 20] | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | R | ✗ | ✗ |
| ARNOLD | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | RT | ✓ | ✓ |

Among works that provide continuous object state change simulation, most do not focus on manipulating object states in a precise and fine-grained manner. For example, VRKitchen [23] defines task goals in a discrete manner even though the underlying object states are continuous. Softgym [46] is an object manipulation benchmark that provides a realistic simulation of deformable objects; however, it lacks diversity among objects and scenes.

By contrast, ARNOLD provides a wide variety of scenes and objects. And ARNOLD simulates continuous states for articulated objects and simulates fluids at the particle level. We control the robots with 7-DOF continuous control and friction-based grasping powered by a state-of-the-art physics engine (PhysX 5.0). Whereas most of the environments [29, 76] optimize for speed, we optimize for the realism of the rendering. And ARNOLD also equips a remarkable rendering speed at 185 FPS (37 FPS with five cameras).

Language Conditioned Manipulation. Relating human language to robot actions has been of recent interest [51, 75, 90, 39, 12, 33, 32, 63]. However, the environments in these efforts either lack realistic physics [73, 71] or do not have realistic scenes [72, 58] where the surroundings of the agent will constrain its motion, and different scene objects might occlude the agent’s viewpoint. Additionally, systems like [33, 4, 14, 52] are application-based, lacking a systematic benchmark for language-conditioned manipulation. Most importantly, prior work aims to ground human language to static object properties, such as colors and shapes. By contrast, ARNOLD provides instructions for continuous object states. We compare between ARNOLD and other related benchmarks in Tab. 1.

Continuous State Understanding. Some recent research tries to predict object states [48, 61]. However, the

object states are discrete rather than continuous. More recently, researchers [83, 13, 82, 78] tried to predict object states continuously. However, they do not address the manipulation of objects from arbitrary starting states to the desired states. Moreover, they do not model the language grounding process. Most recently, Ma *et al.* [53] propose a method to perform precise object state manipulations, but their approach does not perform language grounding and only a small subset of our tasks are covered with their simple motion primitives.

Compared with prior works, we provide more diverse goal states to cover the continuous state space instead of learning only binary goal states. This allows models to understand the continuous state space. On the other hand, we also propose evaluation of generalization in terms of continuous state understanding (see Sec. 3.4). This evaluates how the model leverages its understanding of the state space to generalize within a continuous spectrum, which is rarely studied in prior works. Though more goal states can be added with our continuous simulation, we leave ARNOLD at the current scale since more states will lead to higher costs of data generation due to the compositions of object/scene/state.

3. The ARNOLD Benchmark

The ARNOLD benchmark is motivated by the abilities that an intelligent manipulator agent should possess, including (1) the ability to comprehend and ground human instructions to precise world states, (2) the capacity to acquire policies for generating accurate actions and plans toward precisely defined goal states, and (3) the feasibility of transferring such abilities to the real world. Therefore, in ARNOLD we focus on language-conditioned manipulation driven by continuous

goal states situated in diverse photo-realistic and physically-realistic 3D scenes.

3.1. Simulation Environment

Simulation Platform. ARNOLD is built on NVIDIA’s Isaac Sim [56], a robotic simulation application that provides photo-realistic and physically-accurate simulations for robotics research and development. In ARNOLD, the photo-realistic rendering is powered by GPU-enabled ray tracing, and the physics simulation is based on PhysX 5.0. Fig. 1 and Fig. 2 provide examples of simulation and rendering.

Physical Simulation. To ensure physically-realistic simulation, we assign physics parameters to objects, including weight and friction for rigid-body objects, and cohesion, surface tension, and viscosity for fluids. These parameters are selected as in prior work [60] and are adjusted by human operator feedback. Fluids are simulated using the GPU-accelerated position-based-dynamics (PBD) method [55] through NVIDIA’s Omniverse platform. Depending on the rendering speed, we perform an optional surface construction process using marching cubes [49] to achieve the final fluid rendering effect.

Scene Configuration. There are 40 distinct objects and 20 diverse scenes in ARNOLD. The scenes are curated from [18], a large-scale synthetic dataset of indoor scenes. This endows ARNOLD with professionally designed layouts and high-quality 3D models. In addition to objects provided by Isaac Sim, we collected objects from open-source datasets [41, 85]. We modified object meshes to enhance visual realism, *e.g.*, by modifying materials and adding top covers to cabinets and drawers. For more stable physics-based grasping, we performed convex decomposition to create precise collision proxies for each object. More details are found in Appendix A of [27].

Robot. We use a 7-DoF Franka Emika Panda manipulator with a parallel gripper in ARNOLD for task execution. The agent has direct control over its seven joints and its gripper. We represent end-effector actions with three spatial coordinates for translation and quaternion for rotation, as it is more tractable for policy learning [47]. We utilize the built-in motion planner of Isaac Sim to transform the end-effector action back to the space of robot joints for execution. Currently, our tasks do not involve navigation, *i.e.*, the robot base remains fixed during task execution.

Visual Input. In ARNOLD, we use five cameras around the robot for visual input. As shown in Fig. 2, the cameras provide various views, including front, left, robot base, and wrist. While each camera provides RGB-D input at a resolution of 128×128 by default, users can render at arbitrary resolution. Notably, unlike the deterministic rendering in prior works [72, 90], the rendering in ARNOLD is stochastic due to the ray tracing sampling process [70], which makes ARNOLD more realistic and challenging. In ad-

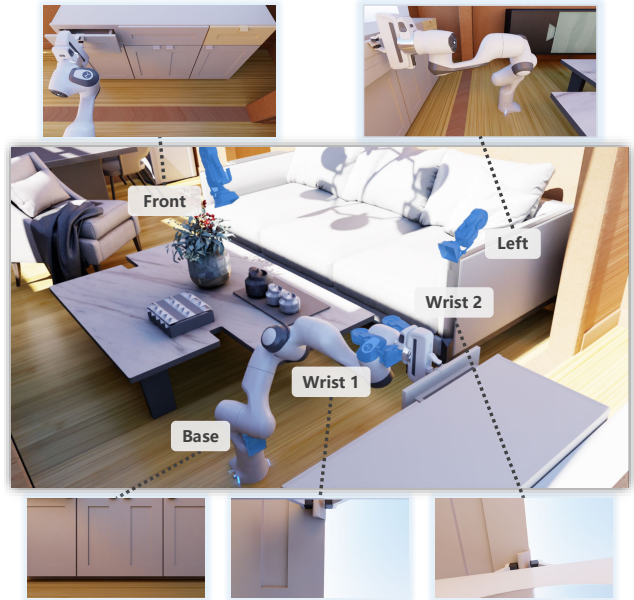


Figure 2. Multi-view robot observation in ARNOLD. The top row shows views from the front and left cameras and the bottom row from the base and two wrist cameras.

dition to the visual observation, other auxiliary observations can be accessed, *e.g.*, camera parameters, robot base pose, and part-level semantic mask. Other Omniverse sensors (*e.g.*, tactile) are excluded here since they are not required by the tasks and models. They are all available if necessary.

3.2. Task Design

We include eight tasks with various goal state variations in ARNOLD. Specifically, we focus on continuous goal states and define success ranges around them wherein robots should maintain object states for 2 seconds to succeed. Tab. 2 provides an overview and Fig. 1 a visualization. More illustrative examples are shown in Appendix B of [27]. Performing these tasks requires capabilities in language grounding, friction-based grasping, continuous state understanding, and robot motion planning. Additional task details follow:

- In PICKUPOBJECT and REORIENTOBJECT, we instruct the robot to manipulate a bottle to achieve different goals. For the former, the initial state of the object is on the ground with goals specifying heights above the ground. For the latter, the initial state of the object is on the ground, oriented horizontally (the state value equals 90°), with goals specifying the angle between the object’s orientation and the upright orientation.
- In the four tasks {OPEN,CLOSE}{DRAWER,CABINET}, the goal value specifies the geometric state of the articulated joint, either in terms of distance (for prismatic joints in DRAWER) or angle (for revolute joints in CABINET). The initial state is any value smaller than the goal for OPEN and larger than the goal for CLOSE.

Table 2. Overview of the 8 tasks in ARNOLD. Each task features 4 goal states specified by human language, one of which is reserved for novel state evaluation. The task is deemed successful when the object state remains in the success range for two seconds. Note that TRANSFERWATER imposes the extra condition that only less than 10% of the original amount of water in the cup can be spilled.

| Task Types | Goal States | Success Ranges |
|----------------|------------------------------|----------------|
| PICKUPOBJECT | 10, 20, 30, 40 (cm) | ± 5 cm |
| REORIENTOBJECT | 0, 45, 135, 180 ($^\circ$) | $\pm 20^\circ$ |
| OPENDRAWER | 25, 50, 75, 100 (%) | $\pm 10\%$ |
| CLOSEDRAWER | 0, 25, 50, 75 (%) | $\pm 10\%$ |
| OPENCABINET | 25, 50, 75, 100 (%) | $\pm 10\%$ |
| CLOSECABINET | 0, 25, 50, 75 (%) | $\pm 10\%$ |
| POURWATER | 25, 50, 75, 100 (%) | $\pm 10\%$ |
| TRANSFERWATER | 20, 40, 60, 80 (%) | $\pm 10\%$ |

- In POURWATER and TRANSFERWATER, the manipulated object is a cup containing water, and the goal specifies the percentage of water to be poured out (POUR) or poured into another cup (TRANSFER). In these two tasks, the goal values are specified as percentages of water relative to the initial amount of water in the cup.

Our task pool covers a variety of manipulation skills and grounding aspects. PICKUPOBJECT and REORIENTOBJECT are selected for the basic skills of moving and rotating objects and the grounding of distances and angles. These abilities are then composed and reinforced in the four tasks {OPEN,CLOSE}{DRAWER,CABINET}, where the goal state is grounded on the state of the manipulated drawer or cabinet joint. Beyond rigid-body objects, fluid manipulation in the two tasks {POUR,TRANSFER}WATER challenges the robots’ ability to manipulate containers and move fluid, grounding goal state values to fluid volumes.

3.3. Data Collection

Demonstration Generation. We designed a motion planner for each task to generate demonstrations. We partitioned each task into sub-task stages for the planner. For each stage, we adopted the RMPflow controller [7] to plan motions toward keypoints. Unlike other approaches to data curation in simulation environments, this keypoint-based motion planner approach affords high sampling efficiency and facilitates imitation learning. While motion planning appeared to be challenging on particular tasks, as demonstrated in [60, 29], we introduced some prior design and practical techniques (details in Appendix B of [27]) to produce satisfactory outcomes. For example, we leveraged spherical linear interpolation (Slerp) to accommodate continuous manipulation in the CABINET and WATER tasks. As a result, our motion planner can efficiently generate demonstrations.

Augmentation With Human Annotations. Despite the strength of motion planners, the diversity of produced demonstrations is highly dependent on the keypoints. To mit-

Table 3. Dataset statistics. (1) *Train*: training data. (2) *Val*: validation data for model selection. (3) *Test*: test data for i.i.d. evaluation. (4) *Object/Scene/State*: *Novel* splits for generalization evaluation. (5) *State**: the *Any State* split for generalization on arbitrary state.

| | <i>Train</i> | <i>Val</i> | <i>Test</i> | <i>Object</i> | <i>Scene</i> | <i>State</i> | <i>State*</i> | Total |
|----------------|--------------|------------|-------------|---------------|--------------|--------------|---------------|--------|
| PICKUPOBJECT | 623 | 134 | 134 | 275 | 221 | 294 | 134 | 1,815 |
| REORIENTOBJECT | 355 | 76 | 77 | 114 | 82 | 210 | 77 | 991 |
| OPENDRAWER | 554 | 119 | 119 | 155 | 255 | 348 | 119 | 1,669 |
| CLOSEDRAWER | 671 | 147 | 148 | 251 | 81 | 530 | 148 | 1,976 |
| OPENCABINET | 319 | 69 | 69 | 81 | 181 | 241 | 69 | 1,029 |
| CLOSECABINET | 478 | 103 | 103 | 55 | 157 | 72 | 103 | 1,071 |
| POURWATER | 312 | 67 | 67 | 96 | 87 | 186 | 67 | 882 |
| TRANSFERWATER | 259 | 56 | 56 | 51 | 50 | 119 | 56 | 647 |
| Total | 3,571 | 771 | 773 | 1,078 | 1,114 | 2,000 | 773 | 10,080 |

igate this problem, we collected about 2k human annotations of task configurations (*e.g.*, object positions), which amount to considerably more diverse and higher quality data. Moreover, we augmented the data with additional relative positions and robot shifts to broaden data variations. Eventually, we curated demonstrations by running inference with ground-truth keypoints and verifying the validity of initial configurations in each execution. In total, we collected 10k valid demonstrations for the ARNOLD benchmark (as in Tab. 3), with each demonstration containing 4–6 keyframes.

Language Instructions. For each demonstration, we sampled a template-based language instruction with our language generation engine. We designed several instruction templates with blanks for each task, and each template can be lexicalized with various phrase candidates. For example, the template “*pull the [position] [object] [percentage] open*” may be lexicalized into “*pull the top drawer 50% open*”. In addition to the representation by explicit numbers, we also prepared a candidate pool of equivalent phrases (*e.g.*, “*fifty percent*”, “*half*”, “*two quarters*”) for random replacement. Note that the instruction does not specify the initial state, so the agent must recognize the current state from the observation. We present template examples in Appendix C of [27].

3.4. Benchmark

Data Split. Evaluating and improving the generalization abilities of robots is a major focus of ARNOLD. To this end, we randomly split the objects, scenes, and goal states into seen and unseen subsets, respectively. We then created the *Normal* split by gathering data with seen objects, scenes, and states. The split was further shuffled and divided into *Train/Val/Test* sets proportioned at 70%/15%/15%. Notably, in addition to providing valid initialization configurations, demonstrations for evaluation splits may be used to provide intermediate ground truth for diagnosing model performance (Sec. 4.3). Furthermore, we created the *Generalization* splits *Novel Object/Scene/State* by gathering data with one of the three components (*i.e.*, objects, scenes, and goal states) unseen; *e.g.*, the *Novel Object* split comprises data of unseen

objects, and seen scenes and states.

While the *Novel State* split addresses the generalization of unseen goal states, we expect that grounding on continuous state representations should help the agent to adapt to any arbitrary state within a continuous range. Therefore, we make the *Any State* split with seen objects and scenes, setting the goal states uniformly distributed over a continuous range, *e.g.*, 0%–100%. Such a design resembles universal tasks with arbitrary goal states and facilitates the evaluation of state generalization. Tab. 3 presents the data statistics.

Metrics. A task instance is regarded as a success when the success condition is satisfied continually for 2 seconds. The success condition requires the current state to be within a tolerance threshold from the goal state; *i.e.*, the success range. The tolerances are derived according to human behaviors and are shown in Tab. 2. Note that TRANSFERWATER imposes the extra condition that only 10% or less of the water can be spilled. The execution of evaluation resembles the composition of sub-task stages in the motion planner (details in Appendix B of [27]). To avoid accidental triggering, we check the success condition after the agent completes the final stage. For example, in the task “*pour half of the water out of the cup*”, the agent succeeds if 40% ~ 60% of the water remains in the cup for 2 seconds after the agent has reoriented the cup upright. We have adopted success rate as the evaluation metric in the ARNOLD.

4. Experiments

4.1. Experimental Setup

Models. To evaluate the existing language-conditioned robotic manipulation models on ARNOLD, we chose two state-of-the-art models as our primary focus: 6D-CLIPort [90] and PerAct [72].

- 6D-CLIPort takes as input an RGB-D image from the top-down view and predicts end-effector poses for the current object and the target action. Each end-effector pose contains an action translation and a categorical prediction over discretized Euler angles. 6D-CLIPort comprises three branches to process the multi-modal input: Transporter-ResNet [88] for the spatial stream, CLIP visual encoder and language encoder [67] for the semantic stream.
- PerAct takes RGB-D images as input to fuse a 3D voxelized grid. In addition, PerAct also requires the proprioception, including gripper states and the current timestep. The proprioception features are tiled on the voxel grid. Next, the hybrid voxel grid is downsampled and flattened to a sequence. Meanwhile, the language instruction is fed to a language encoder (*e.g.*, CLIP [67]) and then appended to the sequence. PerAct uses Perceiver-IO [34] to resample a compact latent representation from the multi-modal long sequence. After decoding, PerAct finally outputs a Q function over the original voxel grid for the prediction

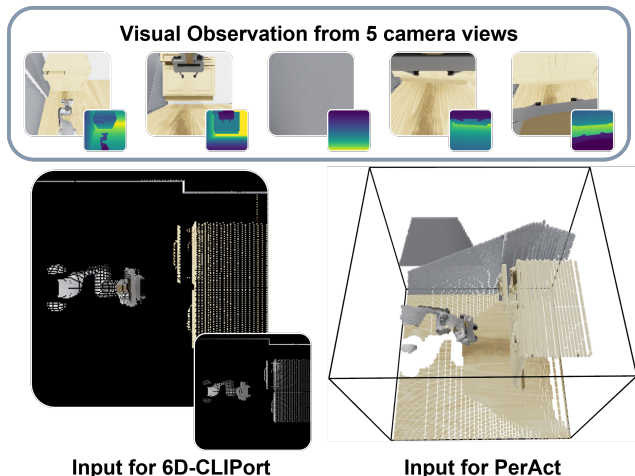


Figure 3. Visualization of the input representations of models. The five cameras are from the front, base, left, and two wrist views. The left camera view is occluded by the wall. Finally, these camera views fuse into an RGB-D image for 6D-CLIPort and a voxel grid for PerAct.

of action translation. Similar to 6D-CLIPort, PerAct also outputs a categorical distribution over discretized Euler angles for the prediction of action rotation. In contrast to the implementation in [72], we discard the heads for predicting gripper and collision. Instead, we add an optional head for state prediction.

Moreover, we considered three model variants of PerAct in our experiments: (1) PerAct without language (PerAct w/o L) for studying the importance of language-grounding, (2) PerAct with additional supervision on state value (PerAct[†]) to show the urgency of state modeling for tasks in ARNOLD, and (3) PerAct trained in the multi-task setting (PerAct MT) given the great potential of multi-task learning shown in [72]. For PerAct[†], we provide additional state supervision by adding an extra output head to regress the normalized state values from the hidden features. We also tried other manipulation models; *e.g.*, BC-Z in Appendix D of [27].

Implementation Details. We obtain the visual representations for the models based on the five rendered RGB-D images as follows: With the camera parameters, we cast the pixels back to 3D and thus derive a point cloud for each view. With these point clouds in 3D scenes, we apply a perception bounding box to make the keypoint learning more tractable. In our setting, the cube spans 126 cm on each axis and the cube center (p_x, p_y, p_z) is 50 cm away from the robot base along the robot’s forward direction. Next, we obtain the visual representations (visualized in Fig. 3) as follows:

- For 6D-CLIPort, with each pixel occupying a size of 0.56 cm, we can map the 126 cm × 126 cm perceived area to a 224 × 224 top-down view image. We project the point clouds with their color and distance information onto this

Table 4. Evaluation results of the models on various tasks and splits, measured by success rate and shown in percentages. The gray figures indicate performances with the first-phase ground truth. For each model, the first row shows the performance on the *Test* set, and the following three rows show those on the *Novel* splits of *Object*, *Scene*, and *State*. The last row indicates the performances on the *Any State* split. Tasks are abbreviated for more space. Average performances on eight tasks are appended to each row. **w/o L**: without language instruction. †: model variants with state modeling. **MT**: multi-task models.

























| | P.OBJECT | R.OBJECT | O.DRAWER | C.DRAWER | O.CABINET | C.CABINET | P.WATER | T.WATER | Average | | | | | | | | | |
|--|----------|----------|----------|----------|-----------|-----------|---------|---------|---------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 6D-CLIPort | 6.72 | 25.37 | 0.00 | 0.00 | 0.00 | 0.00 | 2.70 | 0.00 | 0.00 | 0.00 | 5.83 | 0.00 | 0.00 | 0.00 | 7.14 | 0.84 | 5.13 | |
|  <i>Object</i> | 8.36 | 28.36 | 0.00 | 0.00 | 0.00 | 0.00 | 0.40 | 0.00 | 1.23 | 0.00 | 1.82 | 0.00 | 0.00 | 0.00 | 3.92 | 1.05 | 4.47 | |
|  <i>Scene</i> | 10.41 | 24.43 | 0.00 | 0.00 | 0.00 | 1.57 | 0.00 | 0.00 | 0.55 | 1.27 | 1.27 | 0.00 | 0.00 | 0.00 | 12.00 | 1.46 | 4.98 | |
|  <i>State</i> | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.57 | 0.75 | 1.13 | 0.00 | 0.83 | 0.00 | 2.78 | 0.00 | 0.00 | 16.81 | 0.09 | 2.77 | |
|  <i>Any State</i> | 10.45 | 29.10 | 1.30 | 2.60 | 0.84 | 0.00 | 0.68 | 1.35 | 0.00 | 5.80 | 0.00 | 2.91 | 0.00 | 1.49 | 0.00 | 7.14 | 1.66 | 6.30 |
| PerAct (w/o L) | 25.37 | 33.58 | 14.29 | 7.79 | 17.65 | 36.13 | 47.30 | 52.03 | 8.70 | 34.78 | 7.77 | 10.68 | 14.93 | 11.94 | 5.36 | 14.29 | 17.67 | 25.15 |
|  <i>Object</i> | 29.09 | 26.55 | 8.77 | 3.51 | 3.87 | 20.00 | 24.70 | 32.67 | 0.00 | 0.00 | 1.82 | 7.27 | 16.67 | 29.17 | 9.80 | 19.61 | 11.84 | 17.35 |
|  <i>Scene</i> | 26.70 | 24.89 | 14.63 | 14.63 | 19.61 | 33.33 | 48.15 | 54.32 | 1.10 | 3.87 | 1.91 | 1.27 | 13.79 | 20.69 | 6.00 | 16.00 | 16.49 | 21.13 |
|  <i>State</i> | 0.34 | 0.00 | 0.00 | 0.95 | 9.20 | 8.05 | 1.70 | 2.08 | 0.00 | 2.07 | 1.39 | 1.39 | 1.08 | 2.69 | 5.04 | 8.40 | 2.34 | 3.20 |
|  <i>Any State</i> | 20.15 | 19.40 | 12.99 | 12.99 | 13.45 | 30.25 | 20.95 | 24.32 | 5.80 | 26.09 | 14.56 | 15.53 | 14.93 | 16.42 | 1.79 | 7.14 | 13.08 | 19.02 |
| PerAct | 94.03 | 97.76 | 19.48 | 24.68 | 31.09 | 44.54 | 60.81 | 66.22 | 24.64 | 42.03 | 22.33 | 45.63 | 55.22 | 74.63 | 32.14 | 46.43 | 42.47 | 55.24 |
|  <i>Object</i> | 86.55 | 92.73 | 11.40 | 35.09 | 6.45 | 21.29 | 26.29 | 27.89 | 0.00 | 0.00 | 1.82 | 5.45 | 36.46 | 42.71 | 13.73 | 13.73 | 22.84 | 29.86 |
|  <i>Scene</i> | 72.85 | 84.62 | 17.07 | 31.71 | 20.78 | 31.37 | 66.67 | 64.20 | 0.00 | 4.97 | 5.10 | 19.11 | 33.33 | 51.72 | 22.00 | 36.00 | 29.73 | 40.46 |
|  <i>State</i> | 2.38 | 0.68 | 0.00 | 0.95 | 10.92 | 12.64 | 13.77 | 17.74 | 0.00 | 5.81 | 1.39 | 1.39 | 1.61 | 1.08 | 5.88 | 1.68 | 4.49 | 5.25 |
|  <i>Any State</i> | 47.01 | 50.75 | 7.79 | 19.48 | 21.85 | 30.25 | 18.92 | 25.00 | 5.80 | 21.74 | 3.88 | 15.53 | 14.93 | 25.37 | 10.71 | 14.29 | 16.36 | 25.30 |
| PerAct† | 94.78 | 95.52 | 24.68 | 28.57 | 36.13 | 52.94 | 60.14 | 68.24 | 23.19 | 49.28 | 30.10 | 48.54 | 49.25 | 85.07 | 28.57 | 53.57 | 43.36 | 60.22 |
|  <i>Object</i> | 87.27 | 91.27 | 10.53 | 32.46 | 1.94 | 22.58 | 18.73 | 25.10 | 0.00 | 4.94 | 0.00 | 5.45 | 34.38 | 33.33 | 9.80 | 11.76 | 20.33 | 28.36 |
|  <i>Scene</i> | 69.68 | 84.16 | 13.41 | 37.80 | 25.49 | 38.43 | 60.49 | 67.90 | 0.55 | 4.97 | 6.37 | 19.75 | 29.89 | 63.22 | 26.00 | 24.00 | 28.99 | 42.53 |
|  <i>State</i> | 0.68 | 2.38 | 0.48 | 0.00 | 10.06 | 12.93 | 13.58 | 18.11 | 0.00 | 6.22 | 0.00 | 8.33 | 2.15 | 1.61 | 5.88 | 2.52 | 4.10 | 6.51 |
|  <i>Any State</i> | 48.51 | 47.76 | 14.29 | 14.29 | 21.01 | 33.61 | 23.65 | 28.38 | 4.35 | 24.64 | 6.80 | 13.59 | 26.87 | 31.34 | 14.29 | 19.64 | 19.97 | 26.66 |
| PerAct (MT) | 88.81 | 88.81 | 3.90 | 22.08 | 26.05 | 43.70 | 33.78 | 52.03 | 11.59 | 33.33 | 20.39 | 37.86 | 34.33 | 58.21 | 14.29 | 23.21 | 29.14 | 44.90 |
|  <i>Object</i> | 77.09 | 77.45 | 7.02 | 20.18 | 1.29 | 15.48 | 13.55 | 25.50 | 0.00 | 0.00 | 1.82 | 7.27 | 13.54 | 30.21 | 1.96 | 11.76 | 14.53 | 23.48 |
|  <i>Scene</i> | 68.78 | 77.83 | 10.98 | 41.46 | 13.33 | 25.88 | 29.63 | 55.56 | 0.00 | 4.97 | 5.73 | 9.55 | 19.54 | 43.68 | 4.00 | 16.00 | 19.00 | 34.37 |
|  <i>State</i> | 9.18 | 12.59 | 0.00 | 1.90 | 8.05 | 12.07 | 9.43 | 13.40 | 0.00 | 2.07 | 1.39 | 0.00 | 2.69 | 6.45 | 5.88 | 5.88 | 4.58 | 6.80 |
|  <i>Any State</i> | 36.57 | 45.52 | 5.19 | 10.39 | 15.13 | 19.33 | 16.22 | 21.62 | 1.45 | 8.70 | 0.97 | 3.88 | 8.96 | 14.93 | 3.57 | 17.86 | 11.01 | 17.78 |
| PerAct (MT)† | 90.30 | 92.54 | 14.29 | 20.78 | 25.21 | 47.90 | 33.78 | 56.76 | 20.29 | 39.13 | 19.42 | 37.86 | 26.87 | 64.18 | 17.86 | 30.36 | 31.00 | 48.69 |
|  <i>Object</i> | 81.09 | 85.45 | 7.89 | 24.56 | 3.23 | 22.58 | 14.74 | 28.69 | 0.00 | 4.94 | 1.82 | 5.45 | 9.38 | 20.83 | 3.92 | 19.61 | 15.26 | 26.51 |
|  <i>Scene</i> | 67.87 | 79.64 | 7.32 | 29.27 | 12.94 | 25.49 | 39.51 | 65.43 | 1.10 | 4.42 | 7.01 | 14.01 | 10.34 | 43.68 | 8.00 | 22.00 | 19.26 | 35.49 |
|  <i>State</i> | 2.04 | 3.06 | 0.95 | 2.38 | 9.20 | 18.68 | 6.98 | 11.13 | 0.00 | 3.73 | 2.78 | 11.11 | 6.45 | 9.14 | 1.68 | 4.20 | 3.76 | 7.93 |
|  <i>Any State</i> | 46.27 | 47.01 | 12.99 | 12.99 | 12.61 | 23.53 | 14.86 | 26.35 | 4.35 | 5.80 | 4.85 | 8.74 | 16.42 | 25.37 | 3.57 | 5.36 | 14.49 | 19.39 |

image. Note that the distance represents height rather than depth. If a pixel is related to multiple points (occlusion), the RGB values will be the average of these points and the distance will be the largest height.

- For PerAct, we set the size of voxel grid to be 120^3 , with each voxel covering the 3D context in a 1.05 cm^3 volume. We then aggregate point clouds of all views into the voxel grid, including coordinates, RGB, positional embedding, and occupancy.

Please check <https://github.com/arnold-benchmark/arnold> for detailed implementations of the models and their variants.

Learning. We trained all the models on the *Train* split and performed model selection among ten checkpoints on the *Val* split. Following common practice [90, 72], we used waypoints in demonstrations to facilitate model learning. Specifically, we unify the execution of tasks in ARNOLD into a two-phase procedure: grasping the target object and then manipulating it toward the goal state. We followed the training settings of [90, 72]. Additionally, in multi-task training we sampled each training data by first uniformly sampling the task, then sampling a demonstration of the corresponding task. We set the number of training iterations

to 100k for the single-task and 200k for the multi-task setting. We performed all the training on a single NVIDIA A100 GPU with batch size 8.

Evaluation Execution. The evaluation executor of each task resembles the pipeline of a motion planner yet has a slight difference (see details in Appendix B of [27]). The model’s predictions are converted back to keypoints of sub-task stages for evaluation execution. We make our evaluation strict and appropriate by avoiding shortcuts on the object state during robot motion. For example, the task instructing “pull the cabinet 50% open” will not be considered a success even if the cabinet is held around 50% open for 2 seconds so long as the motion planner has not executed its final action. To eliminate the influence of a first-phase failure (*i.e.*, failing to grasp the target object or object part), we conduct additional evaluations that provide first-phase ground truth.

4.2. Experimental Results

We report experimental results in Tab. 4 (black) and present our findings and analyses below.

Across Models. Comparing the baseline models 6D-CLIPort, PerAct, and PerAct (MT), we found that 6D-

CLIPort fails on most of the tasks. We conjecture that this stems from (1) the information lost in the input representation when compressing complex 3D scenes into a single image and (2) the difficulty in regressing target height values for action translation. By contrast, the voxelized representations in PerAct provide rich 3D contexts that benefit model learning. With such differences, PerAct outperforms 6D-CLIPort significantly. Meanwhile, we observed performance drops for PerAct (MT) compared to PerAct on all tasks. This indicates that it is still difficult to leverage more diverse multi-task data for efficiently learning better task policies in ARNOLD, especially given its challenges in both grounding and manipulation.

Across Tasks. The most challenging tasks in ARNOLD are REORIENTOBJECT and {OPEN,CLOSE}CABINET. REORIENTOBJECT is difficult because it involves estimation of the bottle orientations, and models can often be confused by visually similar states. For example, the action for the goal state of 45° will lead to a goal state of 135° if the bottle orientation is reversed. Manipulating a cabinet proved to be challenging even with privileged information to specify goals [60, 29] as it requires accurate prediction of both interacting position, rotation, and precise continuous motion control. Replacing privileged information with instructions will only make it harder. Furthermore, we observed superior model performance in POURWATER compared to TRANSFERWATER. This is because transferring water requires position alignment between cups to avoid spillage.

On Generalization Splits. In general, we observed performance drops for most models when transferring to *Novel* generalization splits, especially on the *Novel State* split. This reveals that, without proper modeling of continuous states, generalizing the grounding of seen goal states to unseen ones remains challenging. Meanwhile, the performance drop on the *Novel Object* and *Novel Scene* splits varies according to the tasks. For tasks where the objects occupy substantial space (e.g., drawer), the impact of unseen objects is more significant than unseen scenes.

For the *Any State* split, the models’ performances were inferior compared with the *Novel Object/Scene* split and superior to those on the *Novel State* split. As goal states are uniformly sampled from a continuous spectrum, the success ranges of seen goal states are likely to cover a large portion of the spectrum, making the *Any State* generalization interpolations of learned knowledge and skills. This suggests an interesting research question that can be investigated with ARNOLD: How can the task learning model better generalize by interpolating within ranges and extrapolating to out-of-range goal states?

Remarks. The key findings from our experiments with ARNOLD are as follows:

- Current models still struggle with tasks that require complex manipulation skills (e.g., manipulating cabinets). *This*

heightens the demand for better policy learning models to tackle challenging manipulation tasks.

- The low success rate of models on all generalization splits motivates the necessity for (1) *increasingly fine-grained representations for perceptual inputs*, (2) *finer modeling of continuous object states*, and (3) *better alignment between language and robot actions*.
- The *Any state* experiments suggest that *state generalization could potentially be achieved through the interpolation of acquired knowledge and skills*. This promotes approaches with deeper insights into systematic generalization for robot skill adaptation.

4.3. Ablation Studies

Influence of Language. PerAct (w/o L), trained with a single-task scheme, exhibits a considerable performance gap behind PerAct. This indicates the importance of the goal-state information in ARNOLD. Meanwhile, we observe a relatively small gap for REORIENTOBJECT. We believe this is due to (1) the bottleneck of this task lying in the ambiguity of visual perception, as discussed in Sec. 4.2 and (2) the difficulty of grounding angles from current visual observations. On the other hand, the significant performance gap on PICKUPOBJECT shows the effectiveness of language grounding in visually more identifiable concepts such as translation distance.

Importance of State Modeling. The PerAct variants with state supervision (\dagger) were expected to realize a performance gain from explicit state prediction supervision. However, we observed marginal improvements on the *Test* split and limited enhancements on generalization splits for this method. This indicates that such end-to-end state supervision is insufficient for state modeling in ARNOLD and calls for better approaches to representing and modeling the continuous object states in robotic task learning.

With Intermediate Oracle. To better demonstrate how well models understand goal states, we provided models with first-phase ground truth (i.e., grasping positions) and report their performance in Tab. 4 (gray). As expected, we observed significant performance gains with such ground truth. Moreover, directly comparing the scores with the first-phase oracle, we can also observe larger gaps between PerAct and PerAct \dagger , as well as PerAct (MT) and PerAct (MT) \dagger . These results demonstrate that explicit state modeling is indeed beneficial for goal state understanding.

Choice of Language Encoder. We investigated model ablation over the language encoder by switching the default language encoder CLIP [67] in PerAct to T5-base [68]. Due to space constraints, we report model performance only on OPENDRAWER. As shown in Fig. 4, PerAct-T5 outperforms PerAct-CLIP on all the benchmarking splits. This may be due to the inefficacy of the global language representation learned in CLIP in representing fine-grained goal states for

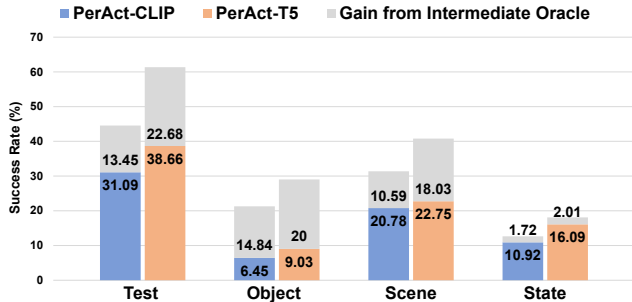


Figure 4. Model ablation results with different language encoders.

precise control. By contrast, T5 is a general-purpose language processing model that offers the ability to maintain detailed information through language modeling and generation. This shows that finer language embeddings may benefit the language grounding of robots to some extent.

4.4. Sim2Real Transfer Experiment

As a realistic simulation environment, one key question to address is: *To what extent can the agents trained in ARNOLD generalize to real-world scenarios?* To this end, we set up a real-world environment for testing the Sim2Real transfer capabilities of agents. Specifically, we used the Franka robot arm to manipulate previously unseen real-world objects with partial point cloud observations captured through a single RGB-D camera from the left view (Fig. 2). We experimented with PerAct, which was trained in ARNOLD to open/close 2 different drawers and pick up 5 different objects. We mitigated the Sim2Real gap as follows:

1. Perception: We used high-fidelity rendering. Due to the imperfect depth sensor in the real world, we sprayed contrast aiding paint onto metallic and transparent areas. For diffuse objects with acceptable depth quality, further domain adaptation would be beneficial.
2. Control: Instead of using the qpos control API, we used an inverse kinematics (IK) controller to perform real-robot actions, which avoids error accumulation.

Additional experimental details can be found in Appendix D of [27]. Throughout our experiments, we observed that models trained in ARNOLD show preliminary Sim2Real transfer capabilities; *i.e.*, reasonable predictions for both picking up objects and manipulating drawers, as shown in Fig. 5. However, the actual robot manipulation continues to struggle because of the Sim2Real gap. For example, when opening a non-plastic drawer in the real world, the robot encounters high friction and is therefore susceptible to prediction errors that lead to inexecutable actions (*e.g.*, exceeding the critical friction angle). With more fine-grained object assets, we believe that the flexible design of the ARNOLD simulator can gradually close this Sim2Real gap by providing more realistic simulations.



Figure 5. Real-world experiments with inference results shown on the upper right. The red dots indicate the predicted positions of the next action for the end-effector.

5. Conclusion

We have presented ARNOLD, a benchmark for language-grounded task learning in realistic 3D interactive environments with diverse scenes, objects, and continuous object states. We devised a systematic benchmark comprising eight challenging language-grounded robot tasks and evaluation splits for robot skill generalization in novel scene, object, and goal-state scenarios. We conducted extensive experiments and analyses to pinpoint the limitations of the current models and identified promising research directions for grounded task learning.

Limitations and Future Work. (1) Despite our focus on realistic simulation, the gap between ARNOLD and real-world scenarios still remains. However, with our flexible design, we can reduce this gap by adding more realistic variations of both scenes and objects [25] to the assets library. (2) Current tasks in ARNOLD do not require much high-level planning knowledge. However, as we found in experiments, short-horizon fine-grained control has not been well solved yet. This defers long-horizon tasks to the future. (3) Despite the diversity we inject in language, the template-based generation inevitably restricts language variations. As a key component for generalization, it is critical to extend ARNOLD with richer language instructions (*e.g.*, by prompting LLMs) in future work. (4) Current imitation learning models rely on prior simplifications (*e.g.*, the two-phase learning) to learn sparse keypoints. Limited by the amount and diversity of expert policies, we need improved methods for modeling continuous object states and learning generalizable policies from scarce data. (5) The realistic and resourceful environment provided by ARNOLD can also facilitate the pursuit of versatile agents with comprehensive capabilities [39, 31]. (6) Scaling up demonstrations is another crucial future step, which can induce stronger capabilities [4, 3].

Acknowledgments

This work is supported in part by the National Key R&D Program of China (2021ZD0150200) and NVIDIA GPU Grant. In addition, we thank Chiao Lu for the help with the environment setup for large-scale experiments.

References

- [1] P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson, N. Sünderhauf, I. Reid, S. Gould, and A. Van Den Hengel, "Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [2] D. Batra, A. X. Chang, S. Chernova, A. J. Davison, J. Deng, V. Koltun, S. Levine, J. Malik, I. Mordatch, R. Mottaghi *et al.*, "Rearrangement: A challenge for embodied AI," *arXiv preprint arXiv:2011.01975*, 2020. 2
- [3] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu *et al.*, "RT-1: Robotics transformer for real-world control at scale," *arXiv preprint arXiv:2212.06817*, 2022. 9
- [4] A. Brohan, Y. Chebotar, C. Finn, K. Hausman, A. Herzog, D. Ho, J. Ibarz, A. Irpan, E. Jang, R. Julian *et al.*, "Do as I can, not as I say: Grounding language in robotic affordances," in *Conference on Robot Learning (CoRL)*, 2022. 3, 9
- [5] A. Chang, A. Dai, T. Funkhouser, M. Halber, M. Niessner, M. Savva, S. Song, A. Zeng, and Y. Zhang, "Matterport3D: Learning from RGB-D data in indoor environments," in *International Conference on 3D Vision (3DV)*, 2017. 2
- [6] C. Chen, C. Schissler, S. Garg, P. Kobernik, A. Clegg, P. Calamia, D. Batra, P. W. Robinson, and K. Grauman, "SoundSpaces 2.0: A simulation platform for visual-acoustic learning," in *Advances in Neural Information Processing Systems Datasets and Benchmarks Track (NeurIPS Datasets and Benchmarks)*, 2022. 2
- [7] C.-A. Cheng, M. Mukadam, J. Issac, S. Birchfield, D. Fox, B. Boots, and N. Ratliff, "RMPflow: A computational graph for automatic motion policy generation," in *Algorithmic Foundations of Robotics XIII: Proceedings of the 13th Workshop on the Algorithmic Foundations of Robotics 13*. Springer, 2020, pp. 441–457. 5
- [8] A. Das, S. Datta, G. Gkioxari, S. Lee, D. Parikh, and D. Batra, "Embodied question answering," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [9] M. Deitke, D. Batra, Y. Bisk, T. Campari, A. X. Chang, D. S. Chaplot, C. Chen, C. P. D'Arpino, K. Ehsani, A. Farhadi *et al.*, "Retrospectives on the embodied AI workshop," *arXiv preprint arXiv:2210.06849*, 2022. 2
- [10] M. Deitke, W. Han, A. Herrasti, A. Kembhavi, E. Kolve, R. Mottaghi, J. Salvador, D. Schwenk, E. VanderBilt, M. Wallingford *et al.*, "RoboTHOR: An open simulation-to-real embodied AI platform," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [11] M. Deitke, E. VanderBilt, A. Herrasti, L. Weihs, J. Salvador, K. Ehsani, W. Han, E. Kolve, A. Farhadi, A. Kembhavi *et al.*, "ProcTHOR: Large-scale embodied AI using procedural generation," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2
- [12] Y. Deng, D. Guo, X. Guo, N. Zhang, H. Liu, and F. Sun, "MQA: Answering the question via robotic manipulation," in *Robotics: Science and Systems (RSS)*, 2020. 3
- [13] Y. Di, R. Zhang, Z. Lou, F. Manhardt, X. Ji, N. Navab, and F. Tombari, "GPV-Pose : Category-level object pose estimation via geometry-guided point-wise voting," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [14] D. Driess, F. Xia, M. S. M. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong, T. Yu, W. Huang, Y. Chebotar, P. Sermanet, D. Duckworth, S. Levine, V. Vanhoucke, K. Hausman, M. Toussaint, K. Greff, A. Zeng, I. Mordatch, and P. Florence, "PaLM-E: An embodied multimodal language model," in *arXiv preprint arXiv:2303.03378*, 2023. 3
- [15] J. Duan, S. Yu, H. L. Tan, H. Zhu, and C. Tan, "A survey of embodied AI: From simulators to research tasks," *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2022. 2
- [16] K. Ehsani, W. Han, A. Herrasti, E. VanderBilt, L. Weihs, E. Kolve, A. Kembhavi, and R. Mottaghi, "ManipulaTHOR: A framework for visual object manipulation," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [17] H. Fu, W. Xu, H. Xue, H. Yang, R. Ye, Y. Huang, Z. Xue, Y. Wang, and C. Lu, "RFUniverse: A physics-based action-centric interactive environment for everyday household tasks," *arXiv preprint arXiv:2202.00199*, 2022. 2
- [18] H. Fu, B. Cai, L. Gao, L.-X. Zhang, J. Wang, C. Li, Q. Zeng, C. Sun, R. Jia, B. Zhao *et al.*, "3D-FRONT: 3D furnished rooms with layouts and semantics," in *International Conference on Computer Vision (ICCV)*, 2021. 2, 4
- [19] C. Gan, J. Schwartz, S. Alter, D. Mrowca, M. Schrimpf, J. Traer, J. De Freitas, J. Kubilius, A. Bhandwaldar, N. Haber *et al.*, "ThreeDWorld: A platform for interactive multi-modal physical simulation," in *Advances in Neural Information Processing Systems Datasets and Benchmarks Track (NeurIPS Datasets and Benchmarks)*, 2021. 2, 3
- [20] C. Gan, S. Zhou, J. Schwartz, S. Alter, A. Bhandwaldar, D. Gutfreund, D. L. Yamins, J. J. DiCarlo, J. McDermott, A. Torralba *et al.*, "The ThreeDWorld transport challenge: A visually guided task-and-motion planning benchmark towards physically realistic embodied AI," in *International Conference on Robotics and Automation (ICRA)*. IEEE, 2022. 3
- [21] Q. Gao, M. Doering, S. Yang, and J. Chai, "Physical causality of action verbs in grounded language understanding," in *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2016. 2

- [22] X. Gao, Q. Gao, R. Gong, K. Lin, G. Thattai, and G. S. Sukhatme, "DialFRED: Dialogue-enabled agents for embodied instruction following," *IEEE Robotics and Automation Letters (RA-L)*, 2022. 2
- [23] X. Gao, R. Gong, T. Shu, X. Xie, S. Wang, and S.-C. Zhu, "VRKitchen: an interactive 3D virtual environment for task-oriented learning," *arXiv preprint arXiv:1903.05757*, 2019. 3
- [24] H. Geng, Z. Li, Y. Geng, J. Chen, H. Dong, and H. Wang, "PartManip: Learning cross-category generalizable part manipulation policy from point cloud observations," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2
- [25] H. Geng, H. Xu, C. Zhao, C. Xu, L. Yi, S. Huang, and H. Wang, "GAPartNet: Cross-category domain-generalizable object perception and manipulation via generalizable and actionable parts," *arXiv preprint arXiv:2211.05272*, 2022. 9
- [26] Y. Geng, B. An, H. Geng, Y. Chen, Y. Yang, and H. Dong, "End-to-end affordance learning for robotic manipulation," *arXiv preprint arXiv:2209.12941*, 2022. 2
- [27] R. Gong, J. Huang, Y. Zhao, H. Geng, X. Gao, Q. Wu, W. Ai, Z. Zhou, D. Terzopoulos, S.-C. Zhu, B. Jia, and S. Huang, "ARNOLD: A benchmark for language-grounded task learning with continuous states in realistic 3D scenes," *arXiv preprint arXiv:2304.04321*, 2023. 4, 5, 6, 7, 9
- [28] J. Gu, D. S. Chaplot, H. Su, and J. Malik, "Multi-skill mobile manipulation for object rearrangement," *International Conference on Learning Representations (ICLR)*, 2022. 2
- [29] J. Gu, F. Xiang, X. Li, Z. Ling, X. Liu, T. Mu, Y. Tang, S. Tao, X. Wei, Y. Yao *et al.*, "ManiSkill2: A unified benchmark for generalizable manipulation skills," in *International Conference on Learning Representations (ICLR)*, 2023. 3, 5, 8
- [30] Y. Hong, Z. Wang, Q. Wu, and S. Gould, "Bridging the gap between learning in discrete and continuous environments for vision-and-language navigation," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [31] J. Huang, W. Y. Zhu, B. Jia, Z. Wang, X. Ma, Q. Li, and S. Huang, "Perceive, ground, reason, and act: A benchmark for general-purpose visual representation," *arXiv preprint arXiv:2211.15402*, 2022. 9
- [32] S. Huang, Z. Wang, P. Li, B. Jia, T. Liu, Y. Zhu, W. Liang, and S.-C. Zhu, "Diffusion-based generation, optimization, and planning in 3D scenes," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3
- [33] W. Huang, F. Xia, T. Xiao, H. Chan, J. Liang, P. Florence, A. Zeng, J. Tompson, I. Mordatch, Y. Chebotar *et al.*, "Inner monologue: Embodied reasoning through planning with language models," in *Conference on Robot Learning (CoRL)*, 2022. 3
- [34] A. Jaegle, S. Borgeaud, J.-B. Alayrac, C. Doersch, C. Ionescu, D. Ding, S. Koppula, D. Zoran, A. Brock, E. Shelhamer *et al.*, "Perceiver IO: A general architecture for structured inputs & outputs," in *International Conference on Learning Representations (ICLR)*, 2022. 6
- [35] U. Jain, L. Weihs, E. Kolve, A. Farhadi, S. Lazebnik, A. Kembhavi, and A. Schwing, "A cordial sync: Going beyond marginal policies for multi-agent embodied tasks," in *European Conference on Computer Vision (ECCV)*, 2020. 2
- [36] U. Jain, L. Weihs, E. Kolve, M. Rastegari, S. Lazebnik, A. Farhadi, A. G. Schwing, and A. Kembhavi, "Two body problem: Collaborative visual task completion," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [37] S. James, Z. Ma, D. R. Arrojo, and A. J. Davison, "RLBench: The robot learning benchmark & learning environment," *IEEE Robotics and Automation Letters (RA-L)*, 2020. 2, 3
- [38] E. Jang, A. Irpan, M. Khansari, D. Kappler, F. Ebert, C. Lynch, S. Levine, and C. Finn, "Bc-z: Zero-shot task generalization with robotic imitation learning," in *Conference on Robot Learning (CoRL)*, 2021. 19
- [39] Y. Jiang, A. Gupta, Z. Zhang, G. Wang, Y. Dou, Y. Chen, L. Fei-Fei, A. Anandkumar, Y. Zhu, and L. Fan, "VIMA: General robot manipulation with multimodal prompts," *arXiv preprint arXiv:2210.03094*, 2022. 3, 9
- [40] A. Kamath, M. Singh, Y. LeCun, G. Synnaeve, I. Misra, and N. Carion, "MDETR-modulated detection for end-to-end multi-modal understanding," in *International Conference on Computer Vision (ICCV)*, 2021. 2
- [41] E. Kolve, R. Mottaghi, W. Han, E. VanderBilt, L. Weihs, A. Herrasti, D. Gordon, Y. Zhu, A. Gupta, and A. Farhadi, "AI2-THOR: An interactive 3D environment for visual ai," *arXiv preprint arXiv:1712.05474*, 2017. 2, 4
- [42] J. Krantz, E. Wijmans, A. Majumdar, D. Batra, and S. Lee, "Beyond the nav-graph: Vision-and-language navigation in continuous environments," in *European Conference on Computer Vision (ECCV)*, 2020. 2
- [43] V. Kumar and E. Todorov, "MuJoCo HAPTIX: A virtual reality system for hand manipulation," in *2015 IEEE-RAS 15th International Conference on Humanoid Robots (Humanoids)*, 2015. 2
- [44] C. Li, F. Xia, R. Martín-Martín, M. Lingelbach, S. Srivastava, B. Shen, K. Vainio, C. Gokmen, G. Dharan, T. Jain *et al.*, "iGibson 2.0: Object-centric simulation for robot learning of everyday household tasks," in *Conference on Robot Learning (CoRL)*, 2021. 2
- [45] C. Li, R. Zhang, J. Wong, C. Gokmen, S. Srivastava, R. Martín-Martín, C. Wang, G. Levine, M. Lingelbach, J. Sun *et al.*, "BEHAVIOR-1K: A benchmark for embodied AI with 1,000 everyday activities and realistic simulation," in *Conference on Robot Learning (CoRL)*, 2022. 2, 3

- [46] X. Lin, Y. Wang, J. Olkin, and D. Held, “SoftGym: Benchmarking deep reinforcement learning for deformable object manipulation,” in *Conference on Robot Learning (CoRL)*, 2020. 2, 3
- [47] M. Liu, X. Li, Z. Ling, Y. Li, and H. Su, “Frame mining: a free lunch for learning robotic manipulation from 3D point clouds,” in *Conference on Robot Learning (CoRL)*, 2022. 4
- [48] Y. Liu, P. Wei, and S.-C. Zhu, “Jointly recognizing object fluents and tasks in egocentric videos,” in *International Conference on Computer Vision (ICCV)*, 2017. 3
- [49] W. E. Lorensen and H. E. Cline, “Marching cubes: A high resolution 3D surface construction algorithm,” *ACM SIGGRAPH Computer Graphics*, vol. 21, no. 4, pp. 163–169, 1987. 4
- [50] J. Lu, D. Batra, D. Parikh, and S. Lee, “ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks,” *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 2
- [51] C. Lynch and P. Sermanet, “Language conditioned imitation learning over unstructured data,” *Robotics: Science and Systems (RSS)*, 2021. 3
- [52] C. Lynch, A. Wahid, J. Tompson, T. Ding, J. Betker, R. Baruch, T. Armstrong, and P. Florence, “Interactive language: Talking to robots in real time,” *arXiv preprint arXiv:2210.06407*, 2022. 3
- [53] L. Ma, J. Meng, S. Liu, W. Chen, J. Xu, and R. Chen, “Sim2Real²: Actively building explicit physics model for precise articulated object manipulation,” in *International Conference on Robotics and Automation (ICRA)*, 2023. 3
- [54] X. Ma, S. Yong, Z. Zheng, Q. Li, Y. Liang, S.-C. Zhu, and S. Huang, “SQA3D: Situated question answering in 3D scenes,” in *International Conference on Learning Representations (ICLR)*, 2022. 2
- [55] M. Macklin and M. Müller, “Position based fluids,” *ACM Transactions on Graphics (TOG)*, vol. 32, no. 4, pp. 1–12, 2013. 4
- [56] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin, A. Allshire, A. Handa *et al.*, “Isaac Gym: High performance GPU-based physics simulation for robot learning,” *Advances in Neural Information Processing Systems Datasets and Benchmarks Track (NeurIPS Datasets and Benchmarks)*, 2021. 4
- [57] O. Mees, L. Hermann, and W. Burgard, “What matters in language conditioned robotic imitation learning over unstructured data,” *IEEE Robotics and Automation Letters (RA-L)*, 2022. 2
- [58] O. Mees, L. Hermann, E. Rosete-Beas, and W. Burgard, “CALVIN: A benchmark for language-conditioned policy learning for long-horizon robot manipulation tasks,” *IEEE Robotics and Automation Letters (RA-L)*, 2022. 2, 3
- [59] M. Mittal, C. Yu, Q. Yu, J. Liu, N. Rudin, D. Hoeller, J. L. Yuan, P. P. Tehrani, R. Singh, Y. Guo *et al.*, “ORBIT: A unified simulation framework for interactive robot learning environments,” *arXiv preprint arXiv:2301.04195*, 2023. 3
- [60] T. Mu, Z. Ling, F. Xiang, D. Yang, X. Li, S. Tao, Z. Huang, Z. Jia, and H. Su, “ManiSkill: Generalizable manipulation skill benchmark with large-scale demonstrations,” *Advances in Neural Information Processing Systems Datasets and Benchmarks Track (NeurIPS Datasets and Benchmarks)*, 2021. 3, 4, 5, 8
- [61] T. Nagarajan and K. Grauman, “Attributes as operators: factorizing unseen attribute-object compositions,” in *European Conference on Computer Vision (ECCV)*, 2018. 3
- [62] —, “Learning affordance landscapes for interaction exploration in 3D environments,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2
- [63] S. Nair, E. Mitchell, K. Chen, S. Savarese, C. Finn *et al.*, “Learning language-conditioned robot behavior from offline data and crowd-sourced annotation,” in *Conference on Robot Learning (CoRL)*, 2022. 3
- [64] A. Padmakumar, J. Thomason, A. Shrivastava, P. Lange, A. Narayan-Chen, S. Gella, R. Piramuthu, G. Tur, and D. Hakkani-Tur, “TEACh: Task-driven embodied agents that chat,” in *AAAI Conference on Artificial Intelligence (AAAI)*, 2022. 2
- [65] A. Pashevich, C. Schmid, and C. Sun, “Episodic transformer for vision-and-language navigation,” in *International Conference on Computer Vision (ICCV)*, 2021. 2
- [66] X. Puig, K. Ra, M. Boben, J. Li, T. Wang, S. Fidler, and A. Torralba, “VirtualHome: Simulating household activities via programs,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [67] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *International Conference on Machine Learning (ICML)*, 2021. 2, 6, 8, 19
- [68] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *Journal of Machine Learning Research (JMLR)*, 2020. 8, 20
- [69] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. Denton, S. K. S. Ghasemipour, B. K. Ayan, S. S. Mahdavi, R. G. Lopes *et al.*, “Photorealistic text-to-image diffusion models with deep language understanding,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2
- [70] P. Shirley and R. K. Morley, *Realistic Ray Tracing*. AK Peters, Ltd., 2008. 4
- [71] M. Shridhar, L. Manuelli, and D. Fox, “CLIPort: What and where pathways for robotic manipulation,” in *Conference on Robot Learning (CoRL)*, 2022. 2, 3, 19

- [72] —, “Perceiver-Actor: A multi-task transformer for robotic manipulation,” *Conference on Robot Learning (CoRL)*, 2022. [2](#), [3](#), [4](#), [6](#), [7](#), [14](#), [19](#)
- [73] M. Shridhar, J. Thomason, D. Gordon, Y. Bisk, W. Han, R. Mottaghi, L. Zettlemoyer, and D. Fox, “ALFRED: A benchmark for interpreting grounded instructions for everyday tasks,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. [2](#), [3](#)
- [74] S. Srivastava, C. Li, M. Lingelbach, R. Martín-Martín, F. Xia, K. E. Vainio, Z. Lian, C. Gokmen, S. Buch, K. Liu *et al.*, “BEHAVIOR: Benchmark for everyday household activities in virtual, interactive, and ecological environments,” in *Conference on Robot Learning (CoRL)*, 2022. [2](#), [3](#)
- [75] S. Stepputtis, J. Campbell, M. Phielipp, S. Lee, C. Baral, and H. Ben Amor, “Language-conditioned imitation learning for robot manipulation tasks,” *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. [3](#)
- [76] A. Szot, A. Clegg, E. Undersander, E. Wijmans, Y. Zhao, J. Turner, N. Maestre, M. Mukadam, D. S. Chaplot, O. Maksymets *et al.*, “Habitat 2.0: Training home assistants to rearrange their habitat,” *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. [2](#), [3](#)
- [77] J. Thomason, M. Murray, M. Cakmak, and L. Zettlemoyer, “Vision-and-dialog navigation,” in *Conference on Robot Learning (CoRL)*, 2020. [2](#)
- [78] W.-C. Tseng, H.-J. Liao, L. Yen-Chen, and M. Sun, “CLANeRF: Category-level articulated neural radiance field,” in *International Conference on Robotics and Automation (ICRA)*, 2022. [3](#)
- [79] W. Wan, H. Geng, Y. Liu, Z. Shan, Y. Yang, L. Yi, and H. Wang, “UniDexGrasp++ : Improving dexterous grasping policy learning via geometry-aware curriculum and iterative generalist-specialist learning,” *arXiv preprint arXiv:2304.00464*, 2023. [2](#)
- [80] X. Wang, Q. Huang, A. Celikyilmaz, J. Gao, D. Shen, Y.-F. Wang, W. Y. Wang, and L. Zhang, “Reinforced cross-modal matching and self-supervised imitation learning for vision-language navigation,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. [2](#)
- [81] Z. Wang, Y. Chen, T. Liu, Y. Zhu, W. Liang, and S. Huang, “HUMANISE: Language-conditioned human motion generation in 3D scenes,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. [2](#)
- [82] F. Wei, R. Chabra, L. Ma, C. Lassner, M. Zollhoefer, S. Rusinkiewicz, C. Sweeney, R. Newcombe, and M. Slavcheva, “Self-supervised neural articulated shape and appearance models,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. [3](#)
- [83] Y. Weng, H. Wang, Q. Zhou, Y. Qin, Y. Duan, Q. Fan, B. Chen, H. Su, and L. J. Guibas, “CAPTRA: Category-level pose tracking for rigid and articulated objects from point clouds,” in *International Conference on Computer Vision (ICCV)*, 2021. [3](#)
- [84] F. Xia, W. B. Shen, C. Li, P. Kasimbeg, M. E. Tchapmi, A. Toshev, R. Martín-Martín, and S. Savarese, “Interactive Gibson benchmark: A benchmark for interactive navigation in cluttered environments,” *IEEE Robotics and Automation Letters (RA-L)*, 2020. [2](#)
- [85] F. Xiang, Y. Qin, K. Mo, Y. Xia, H. Zhu, F. Liu, M. Liu, H. Jiang, Y. Yuan, H. Wang *et al.*, “SAPIEN: A simulated part-based interactive environment,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. [2](#), [4](#)
- [86] E. Xing, A. Gupta, S. Powers, and V. Dean, “KitchenShift: Evaluating zero-shot generalization of imitation-based policy learning under domain shifts,” in *NeurIPS 2021 Workshop on Distribution Shifts: Connecting Methods and Applications*, 2021. [3](#)
- [87] Y. Xu, W. Wan, J. Zhang, H. Liu, Z. Shan, H. Shen, R. Wang, H. Geng, Y. Weng, J. Chen *et al.*, “UniDexGrasp: Universal robotic dexterous grasping via learning diverse proposal generation and goal-conditioned policy,” *arXiv preprint arXiv:2303.00938*, 2023. [2](#)
- [88] A. Zeng, P. Florence, J. Tompson, S. Welker, J. Chien, M. Atarian, T. Armstrong, I. Krasin, D. Duong, V. Sindhwani *et al.*, “Transporter networks: Rearranging the visual world for robotic manipulation,” in *Conference on Robot Learning (CoRL)*, 2021. [3](#), [6](#)
- [89] Y. Zhang and J. Chai, “Hierarchical task learning from language instructions with unified transformers and self-monitoring,” in *Findings of the Association for Computational Linguistics (ACL-IJCNLP Findings)*, 2021. [2](#)
- [90] K. Zheng, X. Chen, O. C. Jenkins, and X. E. Wang, “VLM-bench: A compositional benchmark for vision-and-language manipulation,” *Advances in Neural Information Processing Systems Datasets and Benchmarks Track (NeurIPS Datasets and Benchmarks)*, 2022. [2](#), [3](#), [4](#), [6](#), [7](#), [14](#)

A. Environment

A.1. Assets in USD Format

Universal Scene Description (USD) is a format for 3D scene descriptions. The process of parsing assets into Omniverse involves the transformation of graphical data into USD files. This primarily provides developers with a streamlined method of accessing and retrieving relevant information from the assets, including scene files, articulation bodies, animations, *etc.* By utilizing a USD file format, users can easily access a wide range of assets and subsequently inform Omniverse of the relevant content required for their application, eliminating the need for tedious, manual retrieval of information from each asset individually.

Working with the 3D-Front dataset. The 3D-Front dataset consists of a set of professionally designed synthetic indoor scenes. It features a large number of rooms that are furnished with high-quality textured 3D models. The dataset is composed of three primary components: models, scenes, and textures. And it contains tens of thousands of room layouts with thousands of furnished objects. To parse the 3D-Front dataset into USD format, we apply the following steps:

- Parse the original scene files (JSON) into a data frame containing the mesh and furniture information.
- Use Maya MEL script to load scenes into Autodesk Maya.
- Apply Maya and Omniverse converter to save the scenes in USD format.

Working with articulated bodies. The articulated bodies in ARNOLD mainly come from SAPIEN. We use built-in tools of Omniverse Isaac Sim to parse the original articulated bodies (URDF) into USD format.

A.2. Speed

With five 128×128 cameras, our system runs at 17 fps for liquid simulation and 37 fps for rigid body simulations on an NVIDIA RTX 3090 GPU and AMD 5950X CPU. Currently, Omniverse Isaac Sim only supports serial cameras. We expect a huge performance improvement with the upcoming release of Omniverse Isaac Sim which supports parallel cameras.

A.3. Randomness

Different from previous work [90, 72], the rendering randomness in ARNOLD would result in non-deterministic images. This means two images rendered from the same frame may exhibit a subtle difference. Fig. 8 shows an example where about 78.8% of the pixels are different, together with the distribution of pixel difference.

B. Task Details

We provide illustrative examples of some tasks in Fig. 9. Fig. 10 to 13 illustrate the variations from different aspects

in ARNOLD. In this section, we provide the implementation details of each task. In our notations, we use o to denote visual observation, a for action (*i.e.*, position and rotation with regard to the world). δ is a function to compute the pre-grasp pose.

B.1. PICKUPOBJECT

Motion planner. The motion planner of PICKUPOBJECT consists of four sub-task stages. We depict each stage as follows, beginning with the visual observation and ending with a consequent end effector pose.

1. Observe o_1 . Move the end effector to the pre-grasp pose. Reach $a_1 = \delta(a_2)$.
2. Observe o_2 . Move the end effector to reach the position for grasping. Reach a_2 .
3. Observe o_3 . Close gripper. Reach a_3 .
4. Observe o_4 . Lift the object to the goal height. Reach a_4 .

Learning and evaluation. We extract two observation-action pairs for two-phase learning: o_1, a_2 and o_4, a_4 . During evaluation, the robot executes the action predictions \tilde{a}_2 and \tilde{a}_4 similar to the motion planner:

1. Move to $\delta(\tilde{a}_2)$ for pre-grasp.
2. Move to \tilde{a}_2 and close gripper.
3. Move to \tilde{a}_4 for goal state.

B.2. REORIENTOBJECT

Motion planner. The motion planner of REORIENTOBJECT consists of four sub-task stages. We depict each stage as follows, beginning with the visual observation and ending with a consequent end effector pose.

1. Observe o_1 . Move the end effector to the pre-grasp pose. Reach $a_1 = \delta(a_2)$.
2. Observe o_2 . Move the end effector to reach the position for grasping. Reach a_2 .
3. Observe o_3 . Close gripper. Reach a_3 .
4. Observe o_4 . Reorient the object to goal orientation. Reach a_4 .

Learning and evaluation. We extract two observation-action pairs for two-phase learning: o_1, a_2 and o_4, a_4 . During evaluation, the robot executes the action predictions \tilde{a}_2 and \tilde{a}_4 similar to the motion planner:

1. Move to $\delta(\tilde{a}_2)$ for pre-grasp.
2. Move to \tilde{a}_2 and close gripper.
3. Rotate to \tilde{a}_4 for goal state.

B.3. OPENDRAWER and CLOSEDRAWER

Motion planner. The motion planner for these two tasks consists of four sub-task stages. We depict each stage as follows, beginning with the visual observation and ending with a consequent end effector pose.

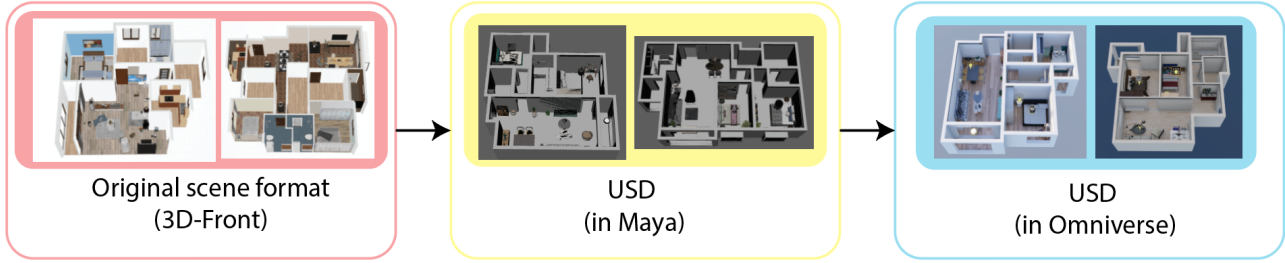


Figure 6. Pipeline of scene parsing. After pre-processing the original 3D-Front scene data, we build an automatic pipeline to load the scene layout into Autodesk Maya with custom designs. Next, we convert the layout file to USD format and deploy it in Omniverse.

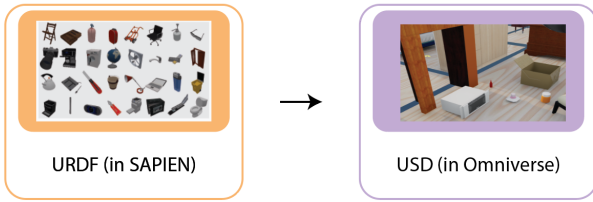
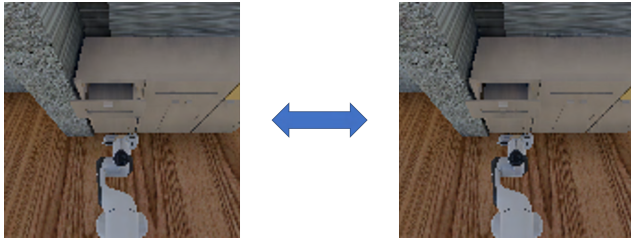


Figure 7. Pipeline of parsing articulated bodies. To convert the original articulated bodies into USD format, we refine the built-in functionalities in Omniverse Isaac Sim and modify the assets manually.

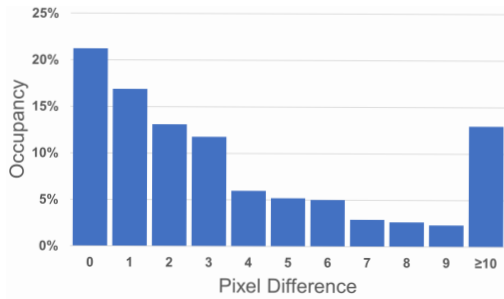
1. Observe o_1 . Move the end effector to the pre-grasp pose. Reach $a_1 = \delta(a_2)$.
2. Observe o_2 . Move the end effector to reach the position for grasping. Reach a_2 .
3. Observe o_3 . Close gripper. Reach a_3 .
4. Observe o_4 . Gradually pull or push the drawer until the goal condition is satisfied. We apply linear interpolation on the action translation to slow down since excessive movement will result in the detachment of gripper. Reach a_4 .

Learning and evaluation. We extract two observation-action pairs for two-phase learning: o_1, a_2 and o_4, a_4 . During evaluation, the robot executes the action predictions \tilde{a}_2 and \tilde{a}_4 similar to the motion planner:

1. Move to $\delta(\tilde{a}_2)$ for pre-grasp.
2. Move to \tilde{a}_2 and close gripper.
3. Gradually interpolate like the motion planner toward \tilde{a}_4 for goal state.



(a) Two rendered images with a subtle difference.



(b) Distribution of pixel difference.

Figure 8. An example that shows the rendering randomness in ARNOLD. (a) Two 128×128 images rendered from the same frame exhibit a subtle difference, with about 78.8% of the pixels being different. (b) The distribution of pixel difference. Here the difference per pixel is measured by the sum of RGB differences. More than 50% of the pixels differ less than a value of 3.

B.4. OPEN CABINET and CLOSE CABINET

Motion planner. The motion planner for these two tasks consists of four sub-task stages. We depict each stage as follows, beginning with the visual observation and ending with a consequent end effector pose.

1. Observe o_1 . Move the end effector to the pre-grasp pose. Reach $a_1 = \delta(a_2)$.
2. Observe o_2 . Move the end effector to reach the position for grasping. Reach a_2 .
3. Observe o_3 . Close gripper. Reach a_3 .
4. Observe o_4 . Gradually pull or push the cabinet until the goal condition is satisfied. We apply linear interpolation on the action translation to slow down, and apply spherical linear interpolation (Slerp) on the action rotation to meet the revolute constraint. Reach a_4 .

Learning and evaluation. We extract two observation-action pairs for two-phase learning: o_1, a_2 and o_4, a_4 . During evaluation, the robot executes the action predictions \tilde{a}_2 and \tilde{a}_4 similar to the motion planner:

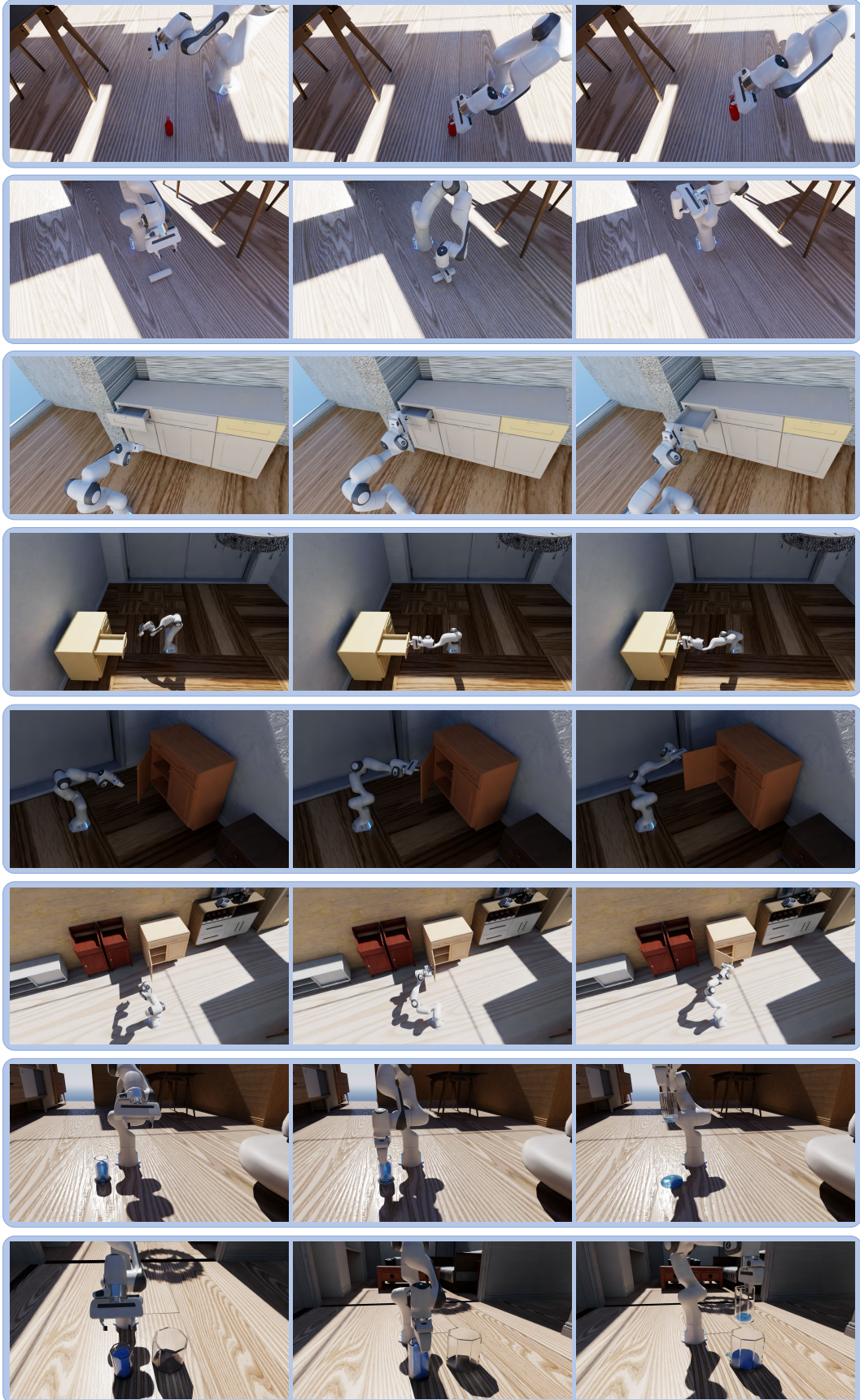


Figure 9. Illustrations of the 8 tasks in ARNOLD.



Figure 10. Scene variations.



Figure 11. Object variations.

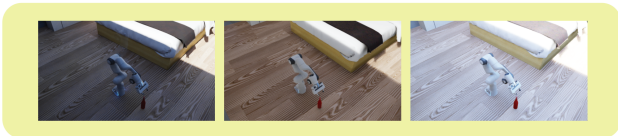


Figure 12. Lighting variations.



Figure 13. Material variations.

1. Move to $\delta(\tilde{a}_2)$ for pre-grasp.
2. Move to \tilde{a}_2 and close gripper.
3. Gradually interpolate like the motion planner toward \tilde{a}_4 for goal state.

B.5. POURWATER and TRANSFERWATER

Motion planner. The motion planner of these two tasks consists of seven sub-task stages. We explain each stage as follows, beginning with the visual observation and ending with a consequent end effector pose:

1. Observe o_1 . Move the end effector to the pre-grasp pose. Reach $a_1 = \delta(a_2)$.
2. Observe o_2 . Move the end effector to reach the position for grasping. Reach a_2 .
3. Observe o_3 . Close gripper. Reach a_3 .
4. Observe o_4 . Lift the cup up to the target height for pouring. Reach a_4 .
5. Observe o_5 . Translate the cup horizontally to the position for pouring. Reach a_5 .
6. Observe o_6 . Gradually tilt the cup to pour water out until the goal condition is satisfied. We apply spherical linear interpolation on the action rotation to slow down and make fluids controllable. Excessive tilting would suddenly empty the water. Reach a_6 .
7. Observe o_7 . Gradually rotate the cup back to an upright pose, similar to the previous stage. Reach a_7 .

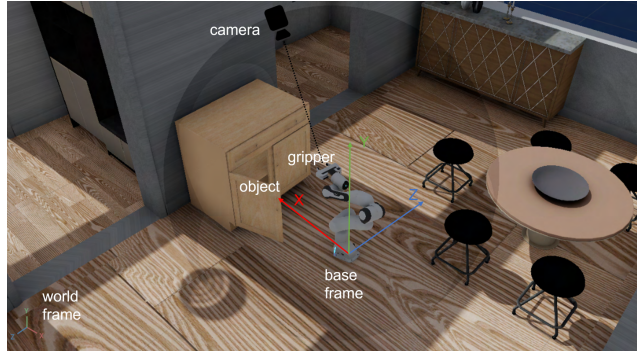


Figure 14. An illustration of the frame and camera for robot teleoperation.

Learning and evaluation. We extract two observation-action pairs for two-phase learning: o_1, a_2 and $o_4, a_{5,6}$, where $a_{5,6}$ combines the translation of a_5 and rotation of a_6 . During evaluation, the robot executes the action predictions $\tilde{a}_2, \tilde{a}_{5,6}$ similar to the motion planner. Note that we maintain the angular velocity to a constant, which ensures reproducibility of the amount of poured water, as well as distinctions among the rotations for different goal states.

1. Move to $\delta(\tilde{a}_2)$ for pre-grasp.
2. Move to \tilde{a}_2 and close gripper.
3. Lift up from \tilde{a}_2 to the height of $\tilde{a}_{5,6}$.
4. Translate horizontally to the position of $\tilde{a}_{5,6}$, with rotation unchanged.
5. Gradually tilt the cup to $\tilde{a}_{5,6}$ by spherical linear interpolation (Slerp).
6. Gradually rotate the cup back to the upright pose.

C. Data Collection

C.1. Human Annotations

We collect human annotations for the positions of robots and interactive assets to ensure reasonable configurations in 3D scenes. For the rationality of such annotations, we ask human operators to teleoperate robots to complete tasks. We will first introduce the settings of robot teleoperation in Appendix C.1.1 and then present the whole pipeline of data collection in Appendix C.1.2.

C.1.1 Robot Teleoperation

Frame definition. Human operators use an Xbox controller to control the robot and the camera for data collection. The input to the controller is in the robot base frame, where the X axis originates from the robot’s base and points to the object, and the Y axis is the upward pointing axis, as displayed in Fig. 14. This controller input is transformed into the world frame for robot motion at each timestep.

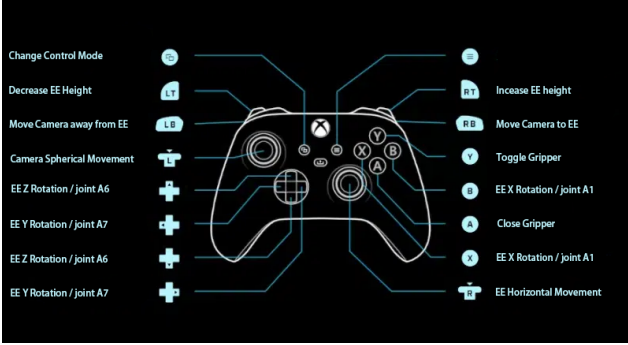


Figure 15. A schematic of the Xbox controller.

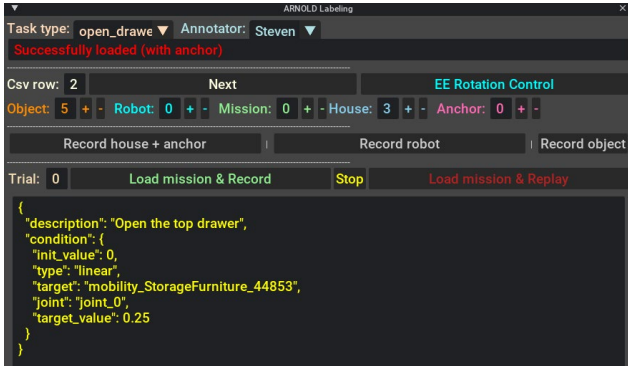


Figure 16. A toy example of the user interface (UI) for collecting human annotations.

Robot control. Human operators can adjust the position and rotation of the robot end effector, as well as toggle the gripper. The Xbox layout is shown in Fig. 15. Specifically, the controller supports two control modes for rotation: joint position control and end effector rotation control. The joint position control mode allows the operator to directly change the positions of joints: A1 (shoulder joint), A6 (fore-arm joint), and A7 (wrist joint). The end effector rotation control mode allows the operator to rotate the end effector around the X, Y, Z axis of the robot base frame while maintaining its position, which is more general for rotation control. Operators can switch between these two modes during the data collection process.

Camera control. Human operators can move their viewing camera freely in the 3D scenes to avoid occlusion. This is accomplished by a spherical camera control (Fig. 14), where the camera can move on a sphere centered at the robot end effector while keeping casting toward the end effector. Note that the radius of the sphere can also be adjusted for a clearer view.

C.1.2 Collection Pipeline

Settings. The collection of human annotations is conducted through a user interface (UI), which is implemented as an extension of Omniverse Isaac Sim, as shown in Fig. 16.

For each task, we enumerate the compositions of objects, scenes, initial states and goal states. Each composition instance is called a mission. Human annotators are supposed to annotate each mission with two configurations of the relative positions between the robot and object. Each configuration is loaded for human control for two trials. Hence, each mission can produce up to four trajectories.

Procedures. The annotator starts annotating a mission by placing the object group (robot and object) at an appropriate anchor position in the scene and clicking the “Record house + anchor” button. Then, the annotator is supposed to adjust the relative position between the robot and object before clicking the “Record robot” and “Record object” buttons. Clicking these two buttons amounts to recording a configuration. For each configuration, the annotator is supposed to control the robot to complete the task after clicking the “Load mission & Record” button. Each click is regarded as a trial. Once the annotator succeeds, a message of “Task Success” will be displayed. Then the annotator can click the “Stop” button and thereby produce a human trajectory. The trajectory replay is supported by clicking the “Load mission & Replay” button. Altogether, each mission is annotated with two configurations and two trials for each configuration. To proceed to the next mission, the annotator can click the “Next” button.

Statistics. We record the human trajectories with the aforementioned Xbox controller at 120Hz and finally collect 2990 trajectories. These trajectories consist of 6.2M frames (14.3 hours) in total and 2073 frames (17.3 seconds) on average. The minimum, median and maximum length are 110 frames (0.9 seconds), 1752 frames (14.6 seconds) and 11336 frames (94.5 seconds), respectively. The distribution is shown in Fig. 17. Although these human trajectories are not directly used for training in ARNOLD, the annotations of configurations and keypoints are referenced by motion planners for demonstration generation. And these human trajectories are of potential use when more powerful algorithms are available.

C.2. Data Augmentation

For richer data variations, we apply augmentation to the robot positions based on collected human annotations. For example, the robot positions may be expanded by shifting 10cm horizontally along four orthogonal directions. After that, we use the motion planner to check if the robot can execute the task successfully with new positions.

C.3. Language instructions

We sample a few instruction templates from our template pool and present them in Tab. 5. For generality and diversity, we construct these templates with placeholders, which can be lexicalized flexibly. In these templates, “[value_object]” holds for the actual object name, e.g., “white bottle”. The

| Tasks | Examples of Delexicalized Templates |
|----------------|--|
| PICKUPOBJECT | <i>Raise</i> [value_object] [value_height] <i>above the ground</i> |
| REORIENTOBJECT | <i>Reorient</i> [value_object] [value_degree] <i>away from the up axis</i> |
| OPENDRAWER | <i>Open the</i> [value_position] [value_object] [value_percent] |
| CLOSEDRAWER | <i>Close the</i> [value_position] [value_object] [value_percent] |
| OPENCABINET | <i>Open the</i> [value_position] [value_object] [value_percent] |
| CLOSECABINET | <i>Close the</i> [value_position] [value_object] [value_percent] |
| POURWATER | <i>Pour</i> [value_percent] <i>water out of</i> [value_object] |
| TRANSFERWATER | <i>Transfer</i> [value_percent] <i>water to</i> [value_object] |

Table 5. A few examples of delexicalized instruction templates for various tasks. During data generation, we first sample a delexicalized instruction template from the template pool. Then, we fill in the placeholders by sampling from candidate pools that consist of equivalent phrases.

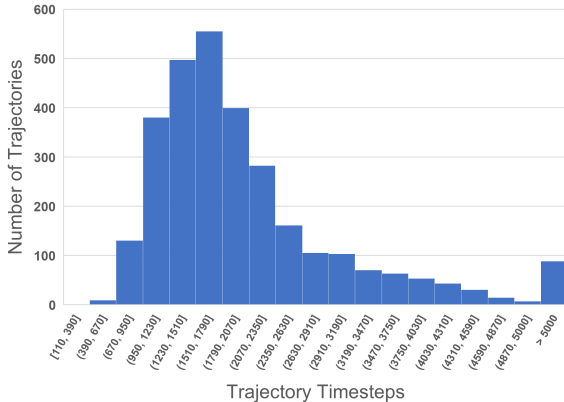


Figure 17. Distribution of human-annotated trajectory length (timesteps). Here 120 timesteps amount to a second.

goal states are specified by “[value_height]” (e.g., “20cm”), “[value_degree]” (e.g., “90 degrees”), and “[value_percent]” (e.g., “40%”). In some tasks where multiple objects exist, the referential words are necessary and specified by “[value_position]”, which indicates the positional information of the target object, e.g., “top left”. The placeholders are lexicalized randomly by sampling from candidate pools which contain equivalent phrases, e.g., “fifty percent”, “half”, “two quarters”.

C.4. Verification

We replay the recorded keypoints of the generated demonstrations in the procedure of evaluation to verify their validity. After removing the failed demonstrations, we obtain a set of 10k demonstrations whose success is ensured achievable.

D. Implementation Details

D.1. Additional Evaluation of BC-Z

BC-Z [38] is a behavior-cloning model that pursues zero-shot task generalization with large-scale data. We follow the implementation of PerAct [72] to adapt the language-

conditioned model (BC-Lang) to our settings. This model takes as input a single-view RGB-D image, with RGB and depth processed by two separate streams. We select the front view for the single-view visual input. The task instruction is processed by a CLIP text encoder [67] to extract a global semantic embedding. With the visual and language input, BC-Lang outputs directly regresses an end effector pose, whose translation and rotation are both continuous values (coordinates and quaternions). There are two variants of backbone: CNN and ViT. We run additional experiments with the two BC-Lang variants and report their performances on the OpenDrawer task (shown in Tab. 6). The results indicate BC-Z cannot handle the tasks in ARNOLD.

D.2. Workflow

Data sampling. We consider three aspects of balance when sampling data for training. (1) To ensure object-level balance, we categorize the demonstrations according to the manipulated object. A uniform sampling over categories is prior to sampling a demonstration within the category. (2) To ensure phase-level balance, we use a biased weight of 1:4 to sample a phase for training based on a selected demonstration. (3) To ensure task-level balance during multi-task training, a uniform sampling over tasks is prior to the task-specific sampling.

Validation. We perform model selection by evaluating each checkpoint instead of tracking the loss on *Val* set since we observe an inconsistency between these two metrics, as mentioned by [71]. We select the best checkpoint for final evaluation.

D.3. Ablation

Without language. We simply skip the operation of appending language embedding to the flattened voxel grid, making the latent representation contain only visual information.

State head. As a naive implementation, we represent all object states through a normalized value in the $[0, 1]$

Table 6. The performances of two BC-Lang variants. We add the performance of PerAct for comparison. The results show that BC-Z, regressing end effector pose with a single-view image as visual input, is much less effective.

| | Test | | Object | | Scene | | State | | Any State | |
|-------------|-------|-------|--------|-------|-------|-------|-------|-------|-----------|-------|
| BC-Lang-CNN | 2.52 | 14.29 | 0.00 | 0.65 | 1.18 | 5.88 | 0.00 | 1.44 | 0.84 | 0.84 |
| BC-Lang-ViT | 5.88 | 19.33 | 0.00 | 0.00 | 0.39 | 2.35 | 0.86 | 4.02 | 0.84 | 5.04 |
| PerAct | 31.09 | 44.54 | 6.45 | 21.29 | 20.78 | 31.37 | 10.92 | 12.64 | 21.85 | 30.25 |

range. For PICKUPOBJECT, the state value is expressed by $\frac{h}{40}$ (initially 0), where h is the height in cm. For REORIENTOBJECT, the state value is expressed by $\frac{\alpha}{180}$ (initially 0.5), where α is the angle (degrees) between object orientation and the upward axis. For the rest of the tasks, the state value is equal to the percentage. The state head predicts both the current state and the goal state, which are optimized via MSE loss. Although learning to regress state values brings moderate improvements, there is still quite a large space for better methods on state modeling.

Language encoder. We adopt T5-base encoder [68] as the alternative to CLIP language encoder. To ensure a fair comparison, we pad the token sequence encoded by T5-base to 77, the same as the token sequence of CLIP. Despite the sophisticated pre-training on large-scale text corpus, T5-base as a language encoder for language-grounded task may be limited by a lack of alignment with vision modality. Moreover, the scarce data in robotics tasks may induce T5-base to suboptimal performances.

D.4. Sim2Real

Equipment and configurations. To set up the real-robot experiment, we adopted the Franka Emika Panda robot arm, the RealSense D435 RGB-D camera, and a random common object in the object category. We placed the robot arm and object similar to the configurations in simulation. Notably, we only used the left camera view, and nothing was tuned particularly. We used the aruco marker to calibrate the D435 camera and Franka API to initialize the robot arm.

Inference. We pre-processed the point cloud due to the imperfect depth camera, *e.g.*, discarding noisy or too far away points. Next, we fetched the PerAct model trained in the simulator to perform inference for the next movement of the robot arm. We utilized Franka API to execute the actions and evaluate similarly to the metrics in ARNOLD.

Results. We summarize the real-world experiment as follows. (1) We experimented with opening/closing two different drawers and picking up five objects. We conducted 19, 8, and 31 trials in various configurations (*e.g.*, object poses, camera poses, instructions *etc.*) for closing drawers, opening drawers, and picking up objects, respectively. And the corresponding numbers of success are 6, 0, and 4. (2) Throughout our experiments, we observe preliminary Sim2Real transfer capabilities, *i.e.*, reasonable predictions for picking up objects and manipulating drawers. Nonetheless, the com-

plex real-world environments, *e.g.*, strict friction and sensory noise, still limit the performance of Sim2Real transfer to a considerable extent.