

PARTINSTRUCT: PART-LEVEL INSTRUCTION FOLLOWING FOR FINE-GRAINED ROBOT MANIPULATION

Anonymous authors

Paper under double-blind review

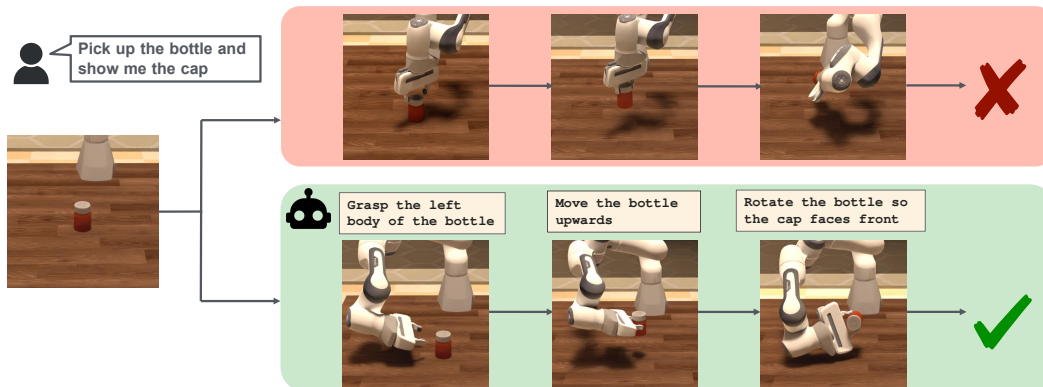


Figure 1: An example fine-grained robot manipulation task in **PartInstruct**. To successfully perform the task described in the instruction (e.g., showing the cap without occluding it), the robot needs to reason about what object parts are relevant, ground the parts to its 3D visual perception, and plan for a sequence of part-level manipulation skills (e.g., the bottom sequence).

ABSTRACT

Fine-grained robot manipulation, such as lifting and rotating a bottle to display the label on the cap, requires robust reasoning about object parts and their relationships with intended tasks. Despite recent advances in training general-purpose robot manipulation policies guided by language instructions, there is a notable lack of large-scale datasets for fine-grained manipulation tasks with part-level instructions and diverse 3D object instances annotated with part-level labels. In this work, we introduce PartInstruct, the first large-scale benchmark for both training and evaluating fine-grained robot manipulation models using part-level instructions. PartInstruct comprises 513 object instances across 14 categories, each annotated with part-level information, and 1302 fine-grained manipulation tasks organized into 16 task classes. Our training set consists of over 10,000 expert demonstrations synthesized in a 3D simulator, where each demonstration is paired with a high-level task instruction, a chain of basic part-based skill instructions, and ground-truth 3D information about the object and its parts. Additionally, we designed a comprehensive test suite to evaluate the generalizability of learned policies across new states, objects, and tasks. We evaluated several state-of-the-art robot manipulation approaches including end-to-end vision-language policy learning and bi-level planning models for robot manipulation on our benchmark. The experimental results reveal that current models struggle to robustly ground part concepts and predict actions in 3D space, and face challenges when manipulating object parts in long-horizon tasks.

1 INTRODUCTION

There has been an increasing interest in training general-purpose vision-language policies for robot manipulation guided by language instructions (Mees et al., 2022; Jiang et al., 2023; Zhang et al., 2023; James et al., 2020; Rahmatizadeh et al., 2018; Zhang et al., 2018), particularly with the recent advances in large generative models (Brohan et al., 2022; Team et al., 2024b). Prior works on language-guided robot manipulation have been mainly focused on high-level manipulation tasks

054 involving simple objects (such as rearranging blocks). However, in the real world, robots often need
 055 to perform fine-grained manipulation of diverse everyday objects, in which the robots need to not
 056 only identify the target object but also understand and interact with specific parts of that object to
 057 perform the intended task as instructed. For instance, to successfully perform the manipulation task
 058 defined in the instruction as shown in Figure 1, the robot needs to identify crucial parts of the object
 059 relevant to the task (such as the label on the cap of the bottle) and reason about a chain of basic
 060 part-based skills that would lead to the desired goal state implied by the instruction, which is to
 061 display the label clearly to the human user without occlusion.

062 Despite the importance of part-level perception and reasoning for robot manipulation, existing robot
 063 manipulation benchmarks on instruction following lack comprehensive integration of part-level se-
 064 mantics in both task instructions and object ground-truth annotations (e.g., James et al., 2020; Jiang
 065 et al., 2023; Mees et al., 2022; Xiang et al., 2020b; Zhang et al., 2023). These benchmarks focus
 066 on object instance-level manipulation tasks but do not include fine-grained, part-level manipulation
 067 tasks like the example in Figure 1. There have been recent benchmarks that evaluate fine-grained,
 068 part-level manipulation tasks, but they either lack language instructions (e.g., Mu et al., 2021; Geng
 069 et al., 2023a) or do not provide training data for policy learning (e.g., Ding et al., 2024).

070 We introduce PartInstruct, the first large-scale benchmark for vision-language policy learning in
 071 fine-grained robot manipulation with part-level semantics. Our core idea is to develop part-level
 072 skills that enable robots to perform complex, fine-grained object manipulation tasks, including those
 073 requiring long-horizon motion plans. To support this, we built *PartGym*, a robot manipulation sim-
 074 ulator for part-level instruction following. Built on PartGym, PartInstruct provides a rich set of
 075 3D object assets with detailed part annotations, along with a large-scale dataset of expert demon-
 076 strations. Additionally, we developed a comprehensive evaluation suite consisting of five test sets,
 077 each corresponding to a different type of generalization test. Together, these tests assess how well a
 078 learned policy performs in unseen scenarios, including new states, objects, and tasks.

079 We evaluated multiple state-of-the-art vision-language policy learning methods designed for
 080 language-guided robot manipulation. We also combined recent learning-based low-level action pol-
 081 icy planning models and VLM-based high-level task planners to create strong bi-level planning
 082 baselines for fine-grained manipulation tasks, which explicitly reasons object parts relevant to a task
 083 and how to interact with them to achieve the final goal. Our experimental results demonstrate that
 084 state-of-the-art methods still struggle with complex fine-grained manipulation tasks.

085 In summary, our main contribution includes (1) the first part-level instruction following benchmark
 086 for both training and evaluating fine-grained robot manipulation models’ capacity for part-level
 087 grounding, reasoning, and planning; (2) a large training dataset with diverse assets and detailed
 088 annotations; (3) a comprehensive evaluation of state-of-the-art vision-language policy learning and
 089 bi-level planning baselines, revealing limitations of current robot manipulation models.

090 2 RELATED WORK

091 2.1 INSTRUCTION FOLLOWING BENCHMARKS FOR TABLE-TOP ROBOT MANIPULATION

092
 093 Early benchmarks like CALVIN (Mees et al., 2022), RLbench (James et al., 2020), and VIMAbench
 094 (Jiang et al., 2023) focused on object-level manipulation but lacked part-level semantics. ManiSkill
 095 (Mu et al., 2021), PartManip (Geng et al., 2023a), and Open6DOR (Ding et al., 2024) introduced
 096 finer object part interactions, with Open6DOR incorporating spatial semantic instructions but rely-
 097 ing on an oracle planner. Recent methods like Composable Part-based Manipulation (CPM) (Liu
 098 et al., 2024), RoboPoint (Yuan et al., 2024), and SAGE (Geng et al., 2023b) emphasize precise part
 099 interactions and semantic grasping, highlighting the need for detailed part-level understanding in
 100 manipulation tasks.
 101

102 2.2 VISION-LANGUAGE POLICIES FOR ROBOT MANIPULATION

103
 104 Vision-language integration in robot manipulation has led to generalist policies like RT-1 (Brohan
 105 et al., 2022), OpenVLA (Kim et al., 2024), and Octo (Team et al., 2024b), which leverage large-scale
 106 vision-language models for versatile instruction execution. Key-pose-based methods, including Per-
 107 Act (Shridhar et al., 2023), Act3D (Gervet et al., 2023), and RVT (Goyal et al., 2023; 2024), focus

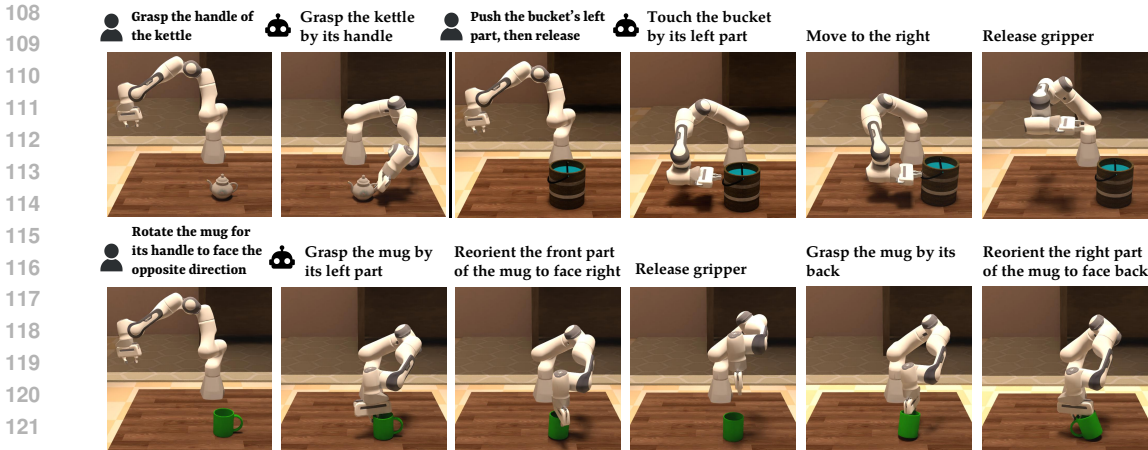


Figure 2: Example tasks and expert demonstrations in the dataset. Each task is defined by a task instruction. Each demonstration is annotated with a chain of base skills and the corresponding skill instructions (the instructions following the task instructions). Specifically, in this figure, the demonstrations for the three tasks have 1, 3, and 5 annotated skill instructions respectively.

on identifying crucial poses to simplify manipulation. Diffusion-based frameworks such as DP (Chi et al., 2023) and DP3 (Ze et al., 2024) employ DDPMs to model multimodal action distributions, generating flexible and expressive robot actions.

2.3 ROBOT PLANNING WITH LLMs AND VLMS.

The integration of LLMs and Vision-Language Models has enhanced robotic planning by improving understanding, reasoning, and task execution. TaPA (Wu et al., 2023) and LLM-Planner (Song et al., 2023) decompose high-level instructions into actionable sub-tasks, while SayCan (Ahn et al., 2022) grounds linguistic instructions in physical affordances to ensure feasible actions. These methods enable robots to interpret and execute multi-step tasks through structured action planning.

3 PARTINSTRUCT BENCHMARK

3.1 PROBLEM SETUP

As shown in Figure 1, a natural language instruction I_{task} describes a part-level instruction following task if it requires that a robot perform a fine-grained manipulation where the robot must interact with a list of object parts in a certain manner to achieve the intended goal g . Critically, the relevant

Table 1: Comparison of PartInstruct with existing tabletop robot manipulation benchmarks based on: the number of distinctive part-level instructions, the number of part labels, the number of fine-grained part-level tasks, availability of training demonstrations, and whether these demonstrations include part-level annotations such as 2D and 3D segmentation masks.

Name	# Part Instruct	# Parts	# Part Tasks	Demo	2D Part Mask	3D Part Mask
CALVIN	6	-	6	✓	✗	✗
RLbench	136	-	64	✗	✗	✗
VIMAbench	-	-	-	✓	✗	✗
LoHoRavens	-	-	-	✓	✗	✗
ManiSkill (SAPIEN)	-	14,068	-	✓	✗	✗
PartManip	-	8,489	1,432	✗	✗	✗
Open6DOR	2,447	-	1,419	✗	✗	✗
PartInstruct (ours)	4,043	4,653	1,302	✓	✓	✓

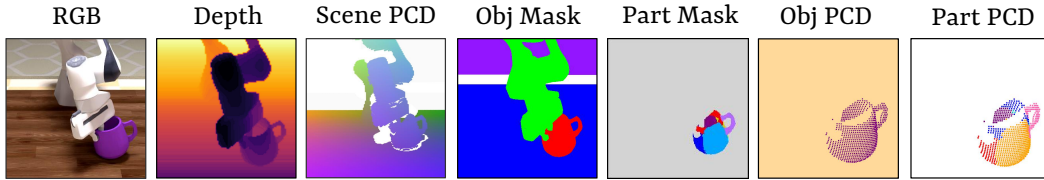


Figure 3: PartGym supports multimodal observations, including RGB images, depth maps, and scene point clouds (PCDs). It also provides object and part annotations, including object segmentations, 2D part segmentation for each object part (part mask), 3D object instance segmentation (obj PCDs), and 3D part segmentations on point clouds (part PCDs) for each object.

Table 2: Definitions of base skills.

Skill	Description
grasp_obj (<i>obj</i> , <i>part</i>)	Robot grasps <i>obj</i> at <i>part</i> .
move_gripper (<i>dir</i> , <i>dis</i> =UNIT, <i>grasping</i> =false)	Robot moves gripper along <i>dir</i> <i>dis</i> .
rotate_obj (<i>obj</i> , <i>part</i> , <i>dir</i>)	Robot rotates <i>obj</i> , such that <i>part</i> is facing <i>dir</i> .
touch_obj (<i>obj</i> , <i>part</i>)	Robot touches <i>obj</i> at <i>part</i> .
release_gripper (<i>obj</i>)	Robot releases the gripper and moves away from <i>obj</i> .

object parts and how the robot needs to interact with them are often not explicitly described in the instructions. Thus the robot must learn to reason about relevant parts and plan how to manipulate them to perform the task successfully. To define g , we first establish a set of goal predicates that specify the states of the object, its parts, the robot’s end effector, and their relationships. For example, **ON** (**obj**, **part**, **surface**) represents physical contact between an object part and a given surface; **FACING** (**obj**, **part**, **dir**) indicates the orientation of an object part from a third-person perspective; and **GRASPING** (**obj**, **part**) denotes a ”grasp” interaction between the object part and the robot’s end effector. Given these goal predicates, each task goal is defined by a set of goal predicates. Examples of tasks are presented in Table 3. In the task illustrated in Figure 1, the goal is represented by the predicate set {**GRASPING** (**bottle**, \sim **cap**), **FACING** (**bottle**, **cap**, **front**), **AT_POSITION** (**bottle**, **INIT_POS+VEC(UP)**)}, where \sim **cap** is any part other than the cap. Note that some tasks consist of multiple phases, where the next phase can only begin after completing the previous one, as the order of interactions is crucial for these tasks. For full task definitions, refer to Appendix A.1.3.

To develop an embodied agent capable of executing tasks defined by g , we hypothesize that it would be beneficial to start with a set of base skills that can be combined to handle a wide range of fine-grained manipulation tasks. In particular, we consider five types of base skills: **grasp_part**, **touch_part**, **rotate_obj**, **move_gripper**, and **release_gripper**. As detailed in Table 2 and Appendix A.1.2, each skill is parameterized by (1) the object part it interacts with and the type of interaction (e.g., touching or grasping), (2) the degree of rotation required for the part, and (3) the distance and direction in which the gripper or object should be moved. This information is summarized in a skill instruction I_{skill} associated with that skill. As illustrated in Figure 2, a task given by an overall task instruction can be decomposed into a sequence of base skill executions. We hypothesize that structuring fine-grained manipulation tasks into sequences of base skills can facilitate the training of hierarchical planning models to compose complex plans with base skills for long-horizon tasks that an end-to-end vision-language policy would struggle with.

3.2 SIMULATION ENVIRONMENT

To train and evaluate language-guided part-level manipulation models, we introduce PartGym, a realistic robot simulator for fine-grained manipulation tasks requiring part-level understanding. PartGym provides (1) rich 3D assets of everyday objects, (2) part-level 3D ground-truth annotations, and (3) a large task set for fine-grained robot manipulation with natural language instructions. We used Pybullet (Coumans & Bai, 2016–2021) as the backbone physics engine to simulate the physical

Table 3: Example task instructions and goal states. Row A corresponds to the task illustrated in Figure 1, while rows B to D correspond to the three tasks shown in Figure 2.

	Task Instruction	Goal States
A	Rotate the <i>part</i> of the <i>object</i> to face <i>direction</i> while lifting it	GRASPING(<i>obj</i>), FACING(<i>part</i> , <i>dir</i>), AT_POSITION(<i>obj</i> , POS_INIT_OBJ+VEC(UP))
B	Grasp the <i>object</i> by the <i>part</i>	GRASPING(<i>gripper</i> , <i>part</i>), ON(<i>obj</i> , <i>table</i>)
C	Move the <i>object</i> to <i>direction</i> by pushing it at the <i>part</i> , then free it	Phase1: TOUCHING(<i>part</i>), AT_POSITION(<i>obj</i> , POS_INIT_OBJ+VEC(<i>dir</i>)) Phase2: MIN_DISTANCE(<i>gripper</i> , <i>obj</i>), GRIPPER_OPEN()
D	Rotate the <i>part</i> of the <i>object</i> to face the opposite direction	FACING(<i>part</i> , ~DIR_INIT(<i>part</i>)), ON(<i>obj</i> , <i>table</i>)

interactions between a robot arm and different objects and their parts. Specifically, the environment includes a 7-DoF Franka Emika Panda robot with a two-finger parallel gripper. Built upon the PartNet Mobility dataset (Xiang et al., 2020a; Mo et al., 2019; Chang et al., 2015), PartGym can simulate manipulation tasks for 14 types of table-top everyday objects annotated with different part labels. In total, there are 513 object instances and 4,653 part labels.

Observations. As shown in Figure 3, we provide multimodal observations for a robot, including RGB images, depth maps, and scene point clouds. Additionally, we provide object and part annotations. Lastly, proprioception robot states like joint states and end-effector poses are also available as part of the observations.

Action Space. The Panda robot takes a 7D action vector at each step. The first 6 dimensions represent the end-effector’s absolute Cartesian pose, parameterized by a 3D coordinate as well as the roll, pitch, and yaw angles. The final dimension controls the gripper’s position.

We provide more details about PartGym in Appendix A.2.

3.3 DATASET

As mentioned in Section 3.1, each fine-grained manipulation task in PartInstruct is accompanied by a natural language description of the overall task, referred to as the **task instruction** I_{task} . Additionally, we provide a sequence of **skill instructions** I_{skill} that specify the part-level manipulation **subgoals** sg to complete the task. It is important to note that skill instructions are provided only during training. For evaluation, models receive only the task instruction as the language input.

3.3.1 TASK CATEGORIES

PartInstruct has 16 task categories, —10 seen and 6 unseen—each requiring specific part-level interactions. Some tasks involve direct part manipulation (e.g. “Hold the part of the object and shift it in a certain direction”), requiring the agent to ground the part in visual data and predict necessary actions. Other task categories require the agent to change the state of a part (e.g. “Rotate the object such that a given part is facing a certain direction”), demanding inference of the final state of the object part and indirect manipulation when needed.

In the 6 test task categories, we have also designed more challenging part-level manipulation tasks. One focus is on long-horizon tasks that require the manipulation of multiple parts in sequence. For instance, “Push the object toward [**direction**] while touching [**part**], lift the object by holding [**part**], then rotate [**part**] to face [**direction**].” Another focus is on tasks that demand more complex reasoning about parts, the environment, and their spatial relationships. For example, consider the task, “Rotate a part of the object on the table so that it points to the opposite direction.” Here, instead of explicitly naming the final state (e.g., a specific direction), the task requires the robot to have additional knowledge about the current direction of a certain part, identify its opposite direction, and manipulate the object so that the part points in that direction.

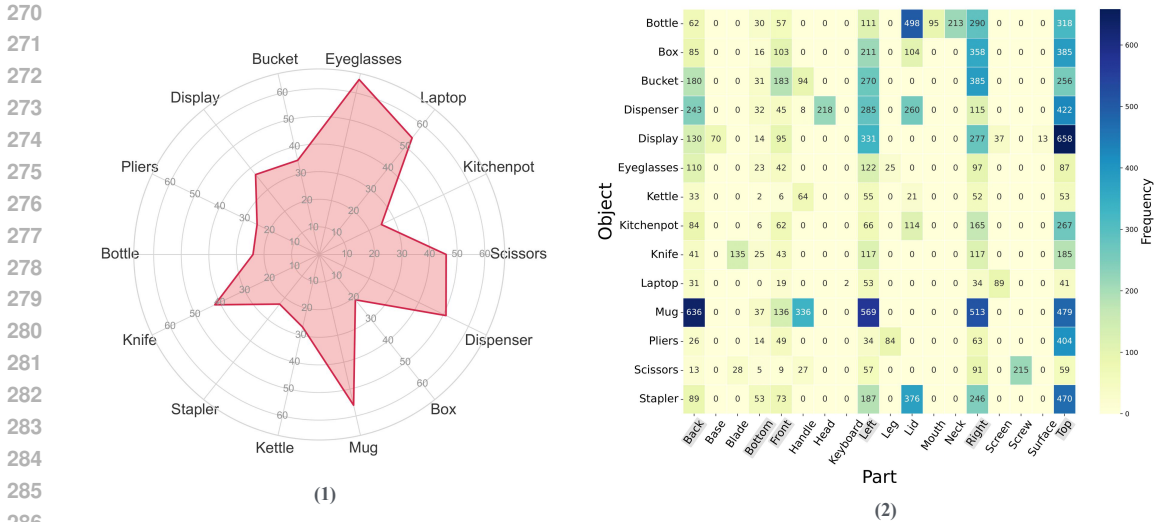


Figure 4: (1) Number of object instances in each object category; (2) Annotated parts grouped by object categories. The horizontal axis stands for different part names and the vertical axis gives different object categories. The value in the heatmap indicates the frequency of each part for an object category in PartInstruct. A darker color shows a higher frequency. Spatial part names are highlighted in light gray to distinguish them from semantic part names.

Table 4: Summary of the five test sets and the type of generalization each one addresses.

Test Set	Type of Generalization
Test 1 (OS)	Novel object positions and rotations
Test 2 (OI)	Novel object instances within the same category
Test 3 (TP)	Novel part combinations within the same task categories
Test 4 (TC)	Novel part-level manipulation task categories
Test 5 (OC)	Novel object categories

3.3.2 TRAJECTORY GENERATION

We leverage a motion planner (use Breyer et al. (2021) for grasping point detection) with oracles to generate a large-scale dataset of demonstrations for vision-language imitation learning. Each demo contains an observation set with different modalities, an expert action trajectory, an overall task instruction that describes the part-level manipulation task in natural language, as well as a chain of skill instructions. See Figure 2 for several example episodes in PartInstruct. Each skill level instruction contains zero or one part that the robot is manipulating with. PartInstruct includes 10,000 demonstrations for training and over 1,800 annotated episodes for testing. Object instance distribution across object categories and the distribution of annotated parts for each object category is shown in Figure 4.

3.3.3 EVALUATION PROTOCOL

Each task defined in Section 3.1 has a binary success criterion. A task is considered complete when every action predicate in its defined predicate set has been satisfied. To systematically evaluate the performance of the learned policy, we designed a five-level evaluation protocol (see Table 4). Each test set evaluates a policy in one type of generalization condition. Specifically, they focus on generalizability over object initial states (OS), novel object instances (OI), novel part combinations in the same task type (TP), novel task categories (TC), and novel object categories (OC). Detailed visualization can be viewed in Appendix A.2.4.

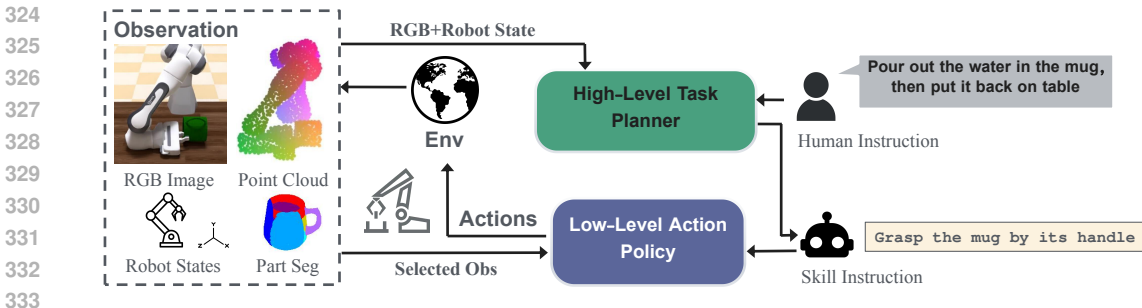


Figure 5: Overview of the bi-level planning framework. The High-Level Task Planner generates a skill instruction as a subgoal for the low-level action policy based on the task instruction and the current observation. Given the subgoal described in the skill instruction, the low-level action policy then generates actions for achieving that subgoal. The high-level task planner updates the skill instruction once every n steps, while the low-level action policy updates the action at every step.

4 EXPERIMENTS

To achieve general-purpose robot manipulation, there have been two common types of approaches: (1) end-to-end policy learning that directly maps observation and instruction to actions (e.g., Zare et al., 2024; Florence et al., 2019; Mandlekar et al., 2020; Rahmatizadeh et al., 2018; Team et al., 2024b; Gervet et al., 2023; Goyal et al., 2024; Chi et al., 2023; Ze et al., 2024) and (2) bi-level planning that first generates high-level plans (typically subgoals), then compute and execute the low-level action plans to achieve the subgoals (e.g., Wu et al., 2023; Song et al., 2023; Ahn et al., 2022; Geng et al., 2023b; Wong et al., 2023). In our benchmark, we evaluate both types of approaches.

4.1 END-TO-END POLICY LEARNING

4.1.1 BASELINES

We evaluate the following state-of-the-art end-to-end robot manipulation policy learning methods:

Octo (Team et al., 2024b) is a transformer-based generalist robot policy pretrained in diverse large-scale robotic episodes.

Act3D (Gervet et al., 2023) is a 3D feature field transformer for multi-task 6-DoF robotic manipulation. Unlike Octo, it employs a key-frame-based approach to complete tasks. During a rollout, the Act3D policy is used to predict the next key pose, which is then executed by commanding the robot to the target key pose using a motion planner.

RVT2 (Goyal et al., 2024) is a multi-task transformer-based 3D manipulation model. Similar to Act3D, it also applies key-frame based manipulation.

Diffusion Policy (DP) (Chi et al., 2023) represents a visuomotor policy as a conditional denoising diffusion process in the action space, which allows it to effectively handle multimodal action distributions and high-dimensional action sequences.

3D Diffusion Policy (DP3) (Ze et al., 2024) combines 3D visual representations with diffusion-based policies, leveraging compact 3D point cloud data for efficient visuomotor policy learning.

Note that the original DP and DP3 models do not support language instruction inputs. To fit the setup of PartInstruct, we modify them to incorporate language inputs. Specifically, we use a pre-trained T5 language encoder to get the language embedding (Raffel et al., 2020). The embedding is then concatenated with other features and used as the observation condition for the denoising process.

We trained the baselines DP, DP3, Act3D, and RVT2 from scratch and fine-tuned the pretrained baseline Octo on our training data. Our hypothesis is that fine-tuning Octo will improve its performance on our benchmark by leveraging its large-scale pretraining on Open X-Embodiment (et al., 2024). The implementation details can be found in Appendix A.4.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

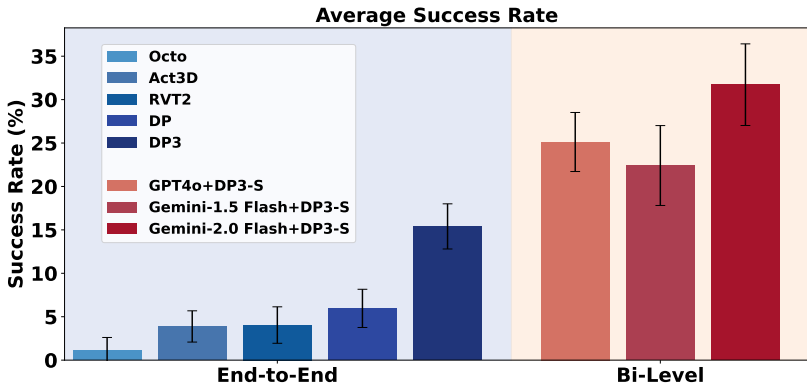


Figure 6: Success Rates of all baselines. The left group represents end-to-end learning policies, while the right group corresponds to bi-level planning models. Error bars denote the standard errors calculated across all evaluation rollouts.

4.1.2 RESULTS

To evaluate each learned policy, we follow the common practice outlined in recent works (Jiang et al., 2023; Chi et al., 2023; Ze et al., 2024). Specifically, we select the top two checkpoints for each baseline and conduct approximately 20 rollouts per object class across all test splits, resulting in over 1,000 rollouts per baseline. We report the *Success Rate* (SR, %) for all end-to-end policy baselines in the left part of Figure 6 and in the top block of Table 5. The low success rate across all baselines suggests that it remains challenging to train an end-to-end generalist policy for fine-grained object manipulation tasks given part-level instructions. They particularly struggle with long-horizon tasks (Test 4) and generalizing to unseen object types (Test 5).

Table 5: Success Rates (%) of baselines across all test sets. The top block includes end-to-end policy learning baselines, and the bottom block includes bi-level planning baselines. The best results are highlighted in blue.

Baselines	OS	OI	TP	TC	OC	All
End-to-End Policy Learning						
Octo	1.82±1.3	0.0	0.91±0.1	0.0	3.33±3.2	1.11±1.5
Act3D	6.25±1.8	5.68±1.7	4.55±1.6	0.0	2.08±2.1	3.88±1.8
RVT2	4.55±2.0	4.55±2.0	6.36±2.3	0.91±0.9	3.33±3.3	4.04±2.1
DP	7.27±1.8	8.64±1.9	8.18±1.8	3.75±2.1	6.67±3.2	5.96±2.2
DP3	23.18±2.8	23.18±2.8	18.18±2.6	7.73±1.8	6.67±3.2	15.40±2.6
Bi-Level Planning						
GPT4o+DP3	33.64±3.2	32.73±3.2	25.91±3.0	10.00±2.0	23.33±5.5	25.12±3.4
Gemini-1.5 Flash+DP3	30.48±4.5	25.45±4.2	27.62±4.4	1.82±1.8	26.67±8.1	22.41±4.6
Gemini-2.0 Flash+DP3	40.58±4.2	34.56±4.1	33.33±4.1	11.90±2.9	38.24±8.3	31.72±4.7

4.2 BI-LEVEL PLANNING

4.2.1 BASELINES

We hypothesize that it would be easier to train action policies with skill instruction annotations compared to directly training a policy for the whole task. Such low-level action policies can then be combined with a high-level planner that generates skill instructions given a task instruction to solve the manipulation task intended by the user. To evaluate the efficacy of bi-level planning on our benchmark, we extend common bi-level planning frameworks (e.g., Geng et al. (2023b)) as shown

Table 6: Performance of low-level action policies when paired with ground-truth high-level plans.

Baselines	OS	OI	TP	TC	OC	All
Octo	3.64	5.50	5.90	0.00	6.67	4.34
Act3D	0.45	2.80	3.33	0.00	0.00	1.32
RVT2	1.82	3.64	1.82	0.00	3.33	1.91
DP	6.47	11.42	16.53	0.06	6.94	8.28
DP3	19.34	13.61	17.89	0.00	12.08	12.58
DP-S	20.00	16.36	25.45	0.00	6.67	13.70
DP3-S	23.64	29.09	23.64	1.82	26.67	20.97

in Figure 5. Specifically, the bi-level planner consists of two modules: (1) a high-level task planner and (2) a low-level action policy. We describe each module below.

High-level Task Planner. We leverage a VLM for high-level task planning. At step t , we prompt the VLM with the task instruction I_{task} to generate the skill instruction for the current step as the subgoal sg_t , i.e., $\pi_{\text{VLM}}(sg_t|o_t, I_{\text{task}})$, where o_t is the observation at step t . We constrain the skill instructions to the space of base skills defined in Table 2 and Appendix A.1.2, which is also specified in the prompt for the VLM. To facilitate decision-making, we also provide additional observations when prompting the VLM, such as RGB images, robot states, etc. See Appendix A.4.3 for the detailed prompt. sg_t will be passed to the low-level action policy for execution and will be updated every n step. Here, n is estimated by the typical length of a skill execution in the training set.

Low-level Action Policy. The low-level action policy is a vision-language policy that generates low-level manipulation actions based on a subgoal and the current observation, i.e., $\pi(a_t|o_t, sg_t)$, where a_t is the action at step t . We can train such policies using the skill instructions annotated for training demonstrations in our dataset. We trained the end-to-end policy learning models evaluated in Section 4.1 on skill instructions as low-level action policies.

Additionally, we hypothesize that an explicit visual understanding of object parts can facilitate part-level instruction grounding. A general vision presentation containing all parts is difficult to obtain since object parts can be overlapped with each other. However, given our benchmark setup as described in Section 3.3.2, the robot interacts with at most one part for the subgoal sg_t defined in each skill instruction, making it possible to give additional vision inputs about the target object part to the low-level action policies. We select the best-performing end-to-end policy learning baselines, DP and DP3, to train the low-level action policies with object part segmentation as part of the input.

For DP, we provide a part segmentation mask as an extra vision input. Given the advanced capability of current general-purpose segmentation models like Segment Anything Model 2 (SAM 2) (Ravi et al., 2024a) in segmenting and tracking object parts, we adopt the approach of Grounded-SAM-2 (Ren et al., 2024). Specifically, given an RGB image and language input, we first utilize a VLM, e.g. Florence-2 (Ravi et al., 2024b) to ground the language onto the target part, then prompt SAM 2 to generate segmentation masks and track the object part in real-time. At each step, we add the obtained part segmentation mask as an extra channel on top of the original RGB, make the input a 4-channel image. We refer to this model as *DP-S*.

For DP3, we use a part point cloud as an additional vision input. Since there has not been a general-purpose object part segmentation model on 3D point cloud (Sun et al., 2024; Sarker et al., 2024), we obtain the 3D part segmentation using a lift-to-3D method. In detail, we first apply the same method in DP to obtain a 2D part mask tracked using SAM2. We then lift the 2D mask into 3D with the depth map and camera intrinsics. To represent a 3D part mask, we append a binary mask channel to the original point cloud observation. We refer to this action policy as *DP3-S*.

The implementation details of bi-level planning baselines can be found in Appendix A.4.3.

4.2.2 RESULTS

We adopt the same evaluation protocol described in Section 4.1.2 for bi-level planning baselines. To evaluate different low-level action policies without considering the effect of high-level task planners, we first pair each low-level action policy with ground-truth skill instructions. As shown in Table 6, *DP3-S* has the highest success rate across all test sets.

Given this result, we then adopt DP3-S as the low-level action policy and pair it with different high-level planners to create bi-level planning baselines. The results are reported in the right part of Figure 6 and the bottom block of Table 5. We can see from the results that the bi-level planning baselines outperform the end-to-end learning in every test set by a large margin. This demonstrates the effectiveness of training a separate low-level action policy for base skills and using VLM as high-level task planner. Among all high-level planning baselines, Gemini-2.0 Flash paired with DP3-S performs the best. However, bi-level planning still struggles with many tasks, particularly when the tasks require longer chains of base skills (e.g., Test 4). In these longer-horizon tasks, there is a higher chance for the high-level task planner to make mistakes. Errors from the low-level action policy are also more likely to be accumulated.

4.3 ABLATION STUDIES

In Section 4.2, we demonstrate that bi-level planning models with low-level action policies informed by part segmentation perform significantly better than state-of-the-art end-to-end policies. To evaluate the effect of each component of the high-level planning models, we conduct the following ablation studies.

Table 7: Impact of high-level task planners on bi-level planning models. We pair each high-level task planner with an oracle motion planner to execute the skill instructions.

Baselines	OS	OI	TP	TC	OC	All
GPT4o	33.12	32.42	34.14	14.42	42.13	31.25
Gemini-1.5 Flash	20.41	19.07	19.36	0.15	29.24	17.65
Gemini-2.0 Flash	27.73	25.94	26.75	0.00	32.70	22.62

4.3.1 EFFECTS OF HIGH-LEVEL PLANNERS

To evaluate different VLMs as high-level planners, we build bi-level planners by combining each VLM with an oracle motion planner—the same one used for training demonstrations—to execute the generated skill instructions. Unlike the full bi-level planning baselines where skill instructions update at fixed intervals, here the oracle planner determines when a subgoal is reached, prompting the high-level planner to issue the next instruction only upon subgoal completion.

Table 7 shows that the bi-level planner using GPT-4o as the high-level task planner improves performance compared to Table 5 when paired with an oracle low-level planner. In contrast, Gemini-based planners perform worse with an oracle, likely because the oracle completes an incorrect skill instruction before Gemini can update it. In such cases, Gemini-based high-level task planners struggle with recovering mistakes made in the previous steps. In contrast, an incorrect skill instruction often leads to the robot arm being stuck in a pose when a learned low-level action policy is trying to perform it.

Table 8: Impact of various vision inputs on low-level action policies. We pair low-level action policies using different vision inputs with ground-truth high-level plans.

Baselines	OS	OI	TP	TC	OC	All
DP	6.47	11.42	16.53	0.06	6.94	8.28
DP-S GT	15.45	20.91	26.36	0.91	13.33	15.39
DP-S SAM2	20.00	16.36	25.45	0.00	6.67	13.70
DP3	19.34	13.61	17.89	0.00	12.08	12.98
DP3-S GT	45.45	36.36	36.36	1.82	40.00	32.00
DP3-S SAM2	23.64	29.09	23.64	1.82	26.67	20.97

4.3.2 EFFECTS OF DIFFERENT VISUAL INPUTS

To examine the impact of different visual representations, particularly 2D and 3D part masks, on policy learning, we conduct another ablation study, where we evaluate the low-level action policies with various visual inputs. Specifically, in addition to *DP-S SAM2* and *DP3-S SAM2*, we also trained low-level action policies using ground-truth mask information, *DP-S GT* and *DP3-S GT*, as well as the vanilla models without any part-level mask, *DP* and *DP3*. The results are summarized

in Table 8. With part segmentations, either 2D or 3D, the low-level action policies can achieve significantly better performance. The performance gap between the policies trained with ground-truth part segmentation and SAM2-based part segmentation also suggests that there improvement in both the VLM’s ability to ground fine-grained parts and in the capacity of state-of-the-art segmentation methods to accurately segment object parts.

4.4 ANALYSIS

4.5 POLICY PERFORMANCE ON PART-LEVEL TASKS

Current vision-language policies perform well on object-level tasks and can follow simple commands such as “grasp” or “touch”. However, they struggle with precise part-level instructions like “touch the left part”, which require fine-grained spatial reasoning. Zero-shot inference using pre-trained generalist policies fails to achieve success (see Appendix A.4.2), likely because these models have not been trained with detailed part-level data. Our dataset and PartGym simulator are valuable because they provide the detailed part annotations and fine-grained tasks for effective training.

4.6 CHALLENGES IN PART-LEVEL INSTRUCTION FOLLOWING

The PartInstruct benchmark shows that part-level instruction following is particularly challenging for state-of-the-art vision-language policies. These models must recognize and track object parts despite variations in appearance—for example, a “lid” may look different on a bottle, pot, stapler, or mug. Moreover, task instructions often do not specify which parts to interact with, forcing the policy to infer the target, while fine-grained tasks demand higher precision and detailed spatial awareness.

Bi-level planning helps address these challenges by decomposing complex tasks into simpler sub-goals that each focus on a single object part. This approach simplifies the training of low-level action policies by reducing the need to track dynamic part-level details and enables the use of pretrained vision-language models for high-level reasoning and planning. Robust visual representations are also crucial for fine-grained manipulation. Our ablation study reveals that 3D representations, such as point clouds, are more effective than 2D images because they provide precise shape and location information. Additionally, explicit object part segmentation—particularly in 3D—significantly boosts performance, with *DP3-S* outperforming *DP3* by approximately 20% (see Table 8).

5 CONCLUSION

In this work, we introduced PartInstruct, a large-scale benchmark designed to advance fine-grained robot manipulation using part-level instructions. By curating a diverse set of objects, tasks, and expert demonstrations, PartInstruct provides a foundation for training and evaluating robot manipulation models that require reasoning about object parts and their relationships with tasks. Our evaluations of state-of-the-art models highlight critical challenges in grounding part concepts and executing long-horizon tasks. With comprehensive experiments, our work provides key insights for future research, highlighting the need for further innovation in perception, reasoning, and planning to enable robots to effectively perform fine-grained, part-aware manipulation.

REFERENCES

- Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- Michel Breyer, Jen Jen Chung, Lionel Ott, Roland Siegwart, and Juan Nieto. Volumetric grasping network: Real-time 6 dof grasp detection in clutter. In *Conference on Robot Learning*, pp. 1602–1611. PMLR, 2021.
- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.

- 594 Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li,
595 Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d
596 model repository. *arXiv preprint arXiv:1512.03012*, 2015.
- 597
- 598 Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shu-
599 ran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *arXiv preprint*
600 *arXiv:2303.04137*, 2023.
- 601 Erwin Coumans and Yunfei Bai. Pybullet, a python module for physics simulation for games,
602 robotics and machine learning. <http://pybullet.org>, 2016–2021.
- 603
- 604 Yufei Ding, Haoran Geng, Chaoyi Xu, Xiaomeng Fang, Jiazhao Zhang, Songlin Wei, Qiyu Dai,
605 Zhizheng Zhang, and He Wang. Open6dor: Benchmarking open-instruction 6-dof object rear-
606 rangement and a vlm-based approach. In *First Vision and Language for Autonomous Driving and*
607 *Robotics Workshop*, 2024.
- 608 Embodiment Collaboration et al. Open x-embodiment: Robotic learning datasets and rt-x models,
609 2024. URL <https://arxiv.org/abs/2310.08864>.
- 610
- 611 Peter Florence, Lucas Manuelli, and Russ Tedrake. Self-supervised correspondence in visuomotor
612 policy learning. *IEEE Robotics and Automation Letters*, 5(2):492–499, 2019.
- 613 Haoran Geng, Ziming Li, Yiran Geng, Jiayi Chen, Hao Dong, and He Wang. Partmanip: Learning
614 cross-category generalizable part manipulation policy from point cloud observations. In *Proceed-*
615 *ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2978–2988,
616 2023a.
- 617 Haoran Geng, Songlin Wei, Congyue Deng, Bokui Shen, He Wang, and Leonidas Guibas. Sage:
618 Bridging semantic and actionable parts for generalizable articulated-object manipulation under
619 language instructions. *arXiv preprint arXiv:2312.01307*, 2023b.
- 620
- 621 Theophile Gervet, Zhou Xian, Nikolaos Gkanatsios, and Katerina Fragkiadaki. Act3d: 3d feature
622 field transformers for multi-task robotic manipulation. In *7th Annual Conference on Robot Learn-*
623 *ing*, 2023.
- 624 Ankit Goyal, Jie Xu, Yijie Guo, Valts Blukis, Yu-Wei Chao, and Dieter Fox. Rvt: Robotic view
625 transformer for 3d object manipulation. In *Conference on Robot Learning*, pp. 694–710. PMLR,
626 2023.
- 627
- 628 Ankit Goyal, Valts Blukis, Jie Xu, Yijie Guo, Yu-Wei Chao, and Dieter Fox. Rvt-2: Learning precise
629 manipulation from few demonstrations. *arXiv preprint arXiv:2406.08545*, 2024.
- 630
- 631 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
632 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
633 770–778, 2016.
- 634 Raisa Islam and Owana Marzia Moushi. Gpt-4o: The cutting-edge advancement in multimodal llm.
635 *Authorea Preprints*, 2024.
- 636
- 637 Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J Davison. Rlbench: The robot
638 learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 5(2):3019–
639 3026, 2020.
- 640 Yunfan Jiang, Agrim Gupta, Zichen Zhang, Guanzhi Wang, Yongqiang Dou, Yanjun Chen, Li Fei-
641 Fei, Anima Anandkumar, Yuke Zhu, and Linxi Fan. Vima: Robot manipulation with multimodal
642 prompts. 2023.
- 643 Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair,
644 Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source
645 vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.
- 646
- 647 Weiyu Liu, Jiayuan Mao, Joy Hsu, Tucker Hermans, Animesh Garg, and Jiajun Wu. Composable
part-based manipulation. *arXiv preprint arXiv:2405.05876*, 2024.

- 648 Ajay Mandlekar, Danfei Xu, Roberto Martín-Martín, Silvio Savarese, and Li Fei-Fei. Learn-
649 ing to generalize across long-horizon tasks from human demonstrations. *arXiv preprint*
650 *arXiv:2003.06085*, 2020.
- 651 Oier Mees, Lukas Hermann, Erick Rosete-Beas, and Wolfram Burgard. Calvin: A benchmark for
652 language-conditioned policy learning for long-horizon robot manipulation tasks. *IEEE Robotics*
653 *and Automation Letters*, 7(3):7327–7334, 2022.
- 654 Kaichun Mo, Shilin Zhu, Angel X. Chang, Li Yi, Subarna Tripathi, Leonidas J. Guibas, and Hao
655 Su. PartNet: A large-scale benchmark for fine-grained and hierarchical part-level 3D object
656 understanding. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,
657 June 2019.
- 658 Tongzhou Mu, Zhan Ling, Fanbo Xiang, Derek Yang, Xuanlin Li, Stone Tao, Zhiao Huang, Zhi-
659 wei Jia, and Hao Su. Maniskill: Generalizable manipulation skill benchmark with large-scale
660 demonstrations. *arXiv preprint arXiv:2107.14483*, 2021.
- 661 Charles R. Qi, Li Yi, Hao Su, and Leonidas J. Guibas. Pointnet++: Deep hierarchical feature
662 learning on point sets in a metric space, 2017. URL <https://arxiv.org/abs/1706.02413>.
- 663 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
664 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
665 models from natural language supervision. In *International conference on machine learning*, pp.
666 8748–8763. PMLR, 2021.
- 667 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
668 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text
669 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 670 Rouhollah Rahmatizadeh, Pooya Abolghasemi, Ladislau Bölöni, and Sergey Levine. Vision-based
671 multi-task manipulation for inexpensive robots using end-to-end learning from demonstration. In
672 *2018 IEEE international conference on robotics and automation (ICRA)*, pp. 3758–3765. IEEE,
673 2018.
- 674 Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham
675 Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Va-
676 sudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Fe-
677 ichtenhofer. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*,
678 2024a. URL <https://arxiv.org/abs/2408.00714>.
- 679 Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham
680 Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Va-
681 sudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Fe-
682 ichtenhofer. Sam 2: Segment anything in images and videos, 2024b. URL <https://arxiv.org/abs/2408.00714>.
- 683 Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang,
684 Yukang Chen, Feng Yan, Zhaoyang Zeng, Hao Zhang, Feng Li, Jie Yang, Hongyang Li, Qing
685 Jiang, and Lei Zhang. Grounded sam: Assembling open-world models for diverse visual tasks,
686 2024. URL <https://arxiv.org/abs/2401.14159>.
- 687 Sushmita Sarker, Prithul Sarker, Gunner Stone, Ryan Gorman, Alireza Tavakkoli, George Bebis,
688 and Javad Sattarvand. A comprehensive overview of deep learning techniques for 3d point cloud
689 classification and semantic segmentation. *Machine Vision and Applications*, 35(4):67, 2024.
- 690 Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Perceiver-actor: A multi-task transformer for
691 robotic manipulation. In *Conference on Robot Learning*, pp. 785–799. PMLR, 2023.
- 692 Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su.
693 Llm-planner: Few-shot grounded planning for embodied agents with large language models. In
694 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2998–3009,
695 2023.

- 702 Yuliang Sun, Xudong Zhang, and Yongwei Miao. A review of point cloud segmentation for under-
703 standing 3d indoor scenes. *Visual Intelligence*, 2(1):14, 2024.
704
- 705 Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer,
706 Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal under-
707 standing across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024a.
- 708 Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep
709 Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist robot
710 policy. *arXiv preprint arXiv:2405.12213*, 2024b.
711
- 712 Lionel Wong, Jiayuan Mao, Pratyusha Sharma, Zachary S Siegel, Jiahai Feng, Noa Korneev,
713 Joshua B Tenenbaum, and Jacob Andreas. Learning adaptive planning representations with natu-
714 ral language guidance. *arXiv preprint arXiv:2312.08566*, 2023.
- 715 Zhenyu Wu, Ziwei Wang, Xiuwei Xu, Jiwen Lu, and Haibin Yan. Embodied task planning with
716 large language models. *arXiv preprint arXiv:2307.01848*, 2023.
717
- 718 Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao
719 Jiang, Yifu Yuan, He Wang, Li Yi, Angel X. Chang, Leonidas J. Guibas, and Hao Su. SAPIEN: A
720 simulated part-based interactive environment. In *The IEEE Conference on Computer Vision and
721 Pattern Recognition (CVPR)*, June 2020a.
- 722 Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao
723 Jiang, Yifu Yuan, He Wang, et al. Sapien: A simulated part-based interactive environment. In
724 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11097–
725 11107, 2020b.
- 726 Wentao Yuan, Jiafei Duan, Valts Blukis, Wilbert Pumacay, Ranjay Krishna, Adithyavairavan Murali,
727 Arsalan Mousavian, and Dieter Fox. Robopoint: A vision-language model for spatial affordance
728 prediction for robotics. *arXiv preprint arXiv:2406.10721*, 2024.
729
- 730 Maryam Zare, Parham M Kebria, Abbas Khosravi, and Saeid Nahavandi. A survey of imitation
731 learning: Algorithms, recent developments, and challenges. *IEEE Transactions on Cybernetics*,
732 2024.
- 733 Yanjie Ze, Gu Zhang, Kangning Zhang, Chenyuan Hu, Muhan Wang, and Huazhe Xu. 3d diffusion
734 policy. *arXiv preprint arXiv:2403.03954*, 2024.
735
- 736 Shengqiang Zhang, Philipp Wicke, Lütfi Kerem Şenel, Luis Figueredo, Abdeldjallil Nacéri,
737 Sami Haddadin, Barbara Plank, and Hinrich Schütze. Lohoravens: A long-horizon language-
738 conditioned benchmark for robotic tabletop manipulation. *arXiv preprint arXiv:2310.12020*,
739 2023.
- 740 Tianhao Zhang, Zoe McCarthy, Owen Jow, Dennis Lee, Xi Chen, Ken Goldberg, and Pieter Abbeel.
741 Deep imitation learning for complex manipulation tasks from virtual reality teleoperation. In
742 *2018 IEEE international conference on robotics and automation (ICRA)*, pp. 5628–5635. IEEE,
743 2018.
744
745
746
747
748
749
750
751
752
753
754
755

A APPENDIX

A.1 PDDL DEFINITIONS

A.1.1 PREDICATE DEFINITIONS

This subsection gives the definition of the basic predicates utilized by the motion planner.

Table 9: Definition of Predicates

Predicate	Description
ON(<i>obj</i> , <i>part</i> , <i>contact</i>)	Whether <i>obj</i> is on the <i>contact</i> .
TOUCHING(<i>obj</i> , <i>part</i>)	Whether the gripper is in contact with <i>obj</i> at <i>part</i> ; <i>part</i> can be empty to indicate a general touch on any parts.
GRASPING(<i>obj</i> , <i>part</i>)	Whether the gripper is carrying <i>obj</i> at <i>part</i> ; <i>part</i> can be empty to indicate a general grasp with any parts.
FACING(<i>obj</i> , <i>part</i> , <i>dir</i>)	Whether <i>part</i> of <i>obj</i> is facing or pointing <i>dir</i> .
AT_POSITION(<i>obj</i> , <i>pos</i>)	Whether <i>obj</i> is at position <i>pos</i> = [<i>x</i> , <i>y</i> , <i>z</i>].

A.1.2 SKILL DEFINITIONS

This subsection shows the detailed definition of the five skills.

Table 10: Definition of Base Skills

Skill	Description	Preconditions	Effects
grasp_obj(<i>obj</i> , <i>part</i>)	Robot grasps <i>obj</i> at <i>part</i> .	ON(<i>table</i> , <i>obj</i>); ~GRASPING(<i>obj</i>); ~TOUCHING(<i>obj</i>)	GRASPING(<i>obj</i> , <i>part</i>)
move_gripper(<i>dir</i> , <i>dis</i> =UNIT, <i>grasping</i> =false)	Robot moves gripper along <i>dir</i> <i>dis</i> .	If <i>grasping</i>==True: GRASPING(<i>obj</i>)	AT_POSITION(<i>gripper</i> , <i>last_gripper_pos</i> + <i>vec</i> (<i>dir</i>) × <i>dis</i>); If <i>grasping</i>==True: GRASPING(<i>obj</i>)
rotate_obj(<i>obj</i> , <i>part</i> , <i>dir</i>)	Robot rotates <i>obj</i> , such that <i>part</i> is facing <i>dir</i> .	GRASPING(<i>obj</i>)	GRASPING(<i>obj</i>); FACING(<i>part</i> , <i>dir</i>)
touch_obj(<i>obj</i> , <i>part</i>)	Robot touches <i>obj</i> at <i>part</i> .	ON(<i>table</i> , <i>obj</i>); ~GRASPING(<i>obj</i>); ~TOUCHING(<i>obj</i>)	TOUCHING(<i>obj</i> , <i>part</i>)
release_gripper(<i>obj</i>)	Robot releases the gripper and moves away from <i>obj</i> .	ON(<i>table</i> , <i>obj</i>); GRASPING(<i>obj</i>) <i>or</i> TOUCHING(<i>obj</i>)	ON(<i>table</i> , <i>obj</i>); ~GRASPING(<i>obj</i>); ~TOUCHING(<i>obj</i>)

A.1.3 TASK DEFINITIONS

This subsection shows the detailed definition of different task types in PartInstruct.

Table 11: Seen Task Instructions and Goal States

Seen (10)		
Order	Example Task Instruction	Goal States
1	Grasp the <i>object</i> by the <i>part</i>	GRASPING(<i>gripper</i> , <i>part</i>), ON(<i>obj</i> , <i>table</i>)
2	Touch the <i>object</i> at the <i>part</i>	TOUCHING(<i>part</i>), ON(<i>obj</i> , <i>table</i>)
3	Hold the <i>part</i> of the <i>object</i> and move it to <i>direction</i>	GRASPING(<i>part</i>), AT.POSITION(<i>obj</i> , POS_INIT.OBJ+VEC(<i>dir</i>))
4	Push the <i>object</i> towards <i>direction</i> by touching <i>part</i>	TOUCHING(<i>part</i>), AT.POSITION(<i>obj</i> , POS_INIT.OBJ+VEC(<i>dir</i>))
5	Slide the <i>object</i> on the table towards <i>direction</i> while keeping hold of <i>part</i> , then release it	<i>Phase1</i> : GRASPING(<i>part</i>), AT.POSITION(<i>obj</i> , POS_INIT.OBJ+VEC(<i>dir</i>)) <i>Phase2</i> : GRIPPER.OPEN, MIN.DISTANCE(<i>gripper</i> , <i>obj</i>)
6	Move the <i>object</i> to <i>direction</i> by pushing it at <i>part</i> , then free it	<i>Phase1</i> : TOUCHING(<i>part</i>), AT.POSITION(<i>obj</i> , POS_INIT.OBJ+VEC(<i>dir</i>)) <i>Phase2</i> : GRIPPER.OPEN, MIN.DISTANCE(<i>gripper</i> , <i>obj</i>)
7	While keeping hold of <i>part</i> , move the <i>object</i> towards <i>direction</i> in the air	GRASPING(<i>part</i>), AT.POSITION(<i>obj</i> , POS_INIT.OBJ+VEC(<i>dir</i>)+VEC(UP))
8	Rotate <i>part</i> of the <i>object</i> to face <i>direction</i> while lifting it	GRASPING(<i>obj</i>), FACING(<i>part</i> , <i>dir</i>), AT.POSITION(<i>obj</i> , POS_INIT.OBJ+VEC(UP))
9	Move the <i>object</i> towards <i>direction</i> after raising it, while keeping hold of <i>part</i> , then put it down	<i>Phase1</i> : GRASPING(<i>part</i>), AT.POSITION(<i>obj</i> , POS_INIT.OBJ+VEC(<i>dir</i>)+VEC(UP)) <i>Phase2</i> : GRASPING(<i>part</i>), AT.POSITION(<i>obj</i> , POS_INIT.OBJ+VEC(<i>dir</i>))
10	Move the <i>object</i> towards <i>direction1</i> in the air, then rotate <i>part</i> to point towards <i>direction2</i>	GRASPING(<i>obj</i>), AT.POSITION(<i>obj</i> , POS_INIT.OBJ+VEC(UP)+VEC(<i>dir1</i>)), FACING(<i>part</i> , <i>dir2</i>)

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

Table 12: Unseen Task Instructions and Goal States

Unseen (6)		
Order	Example Task Instruction	Goal States
11	Rotate <i>part</i> in the air so it points towards <i>direction</i> , then put it down	<i>Phase1:</i> GRASPING(obj), FACING(part, dir), AT_POSITION(obj, POS_INIT.OBJ+VEC(UP)) <i>Phase2:</i> GRASPING(obj), FACING(part, dir), AT_POSITION(obj, POS_INIT.OBJ)
12	Shift the <i>object</i> towards <i>direction1</i> in the air while grasping <i>part1</i> , turn <i>part2</i> to <i>direction2</i> , then set it down	<i>Phase1:</i> GRASPING(part1), FACING(part, dir), AT_POSITION(obj, POS_INIT.OBJ+VEC(dir1)+VEC(UP)) <i>Phase2:</i> GRASPING(part1), FACING(part, dir), AT_POSITION(obj, POS_INIT.OBJ+VEC(dir1))
13	Turn <i>part</i> of the <i>object</i> to point to <i>direction1</i> while keeping it on the table, then push it towards <i>direction2</i>	<i>Phase1:</i> ON(obj, table), FACING(part, dir1); <i>Phase2:</i> ON(obj, table), AT_POSITION(obj, POS_INIT.OBJ+VEC(dir2))
14	While keeping it on the table, push the <i>object</i> towards <i>direction1</i> while touching <i>part1</i> , then rotate <i>part2</i> to face <i>direction2</i>	<i>Phase1:</i> ON(obj, table), TOUCHING(part1), AT_POSITION(obj, POS_INIT.OBJ+VEC(dir1)) <i>Phase2:</i> ON(obj, table), FACING(part2, dir2)
15	Rotate <i>part</i> of the <i>object</i> to face the opposite direction	FACING(part, ~DIR_INIT(part)), ON(obj, table)
16	Push the <i>object</i> to <i>direction1</i> and rotate <i>part</i> to point towards <i>direction2</i> in the air, finally place it down	<i>Phase1:</i> FACING(part, dir2), AT_POSITION(obj, POS_INIT.OBJ+VEC(dir1)+VEC(UP)) <i>Phase2:</i> FACING(part, dir2), AT_POSITION(obj, POS_INIT.OBJ+VEC(dir1))

918 A.2 BENCHMARK

919 A.2.1 OBSERVATION AND ACTION SPACE

920 Table 13 shows the observation and action space available in *PartGym*.

921 Table 13: Observation and Action Space details.

Observation Space	
Static View - RGB	$300 \times 300 \times 3$
Static View - Depth	300×300
Static View - PCD	3×1024
Static View - Semantic	300×300
Static View - Target Part PCD	3×1024
Static View - Target Part Mask	300×300
Wrist View - RGB	$300 \times 300 \times 3$
Wrist View - Depth	300×300
Wrist View - PCD	3×1024
Wrist View - Semantic	300×300
Wrist View - Target Part Mask	300×300
Wrist View - Target Part PCD	3×1024
Proprioceptive state	EE position (3)
	EE orientation (3)
	Joint positions (7)
	Gripper action (1)
Action Space	
Absolute cartesian pose (w.r.t. world frame)	EE position (3)
	EE orientation (3)
	Gripper action (1)

941 A.2.2 TABLE-TOP OBJECTS IN PARTGYM

942 PartGym provides in total 513 object instances across 14 categories. Some objects and their textures

943 are shown in Figure 7.



944 Figure 7: Different object instances in PartGym.

945

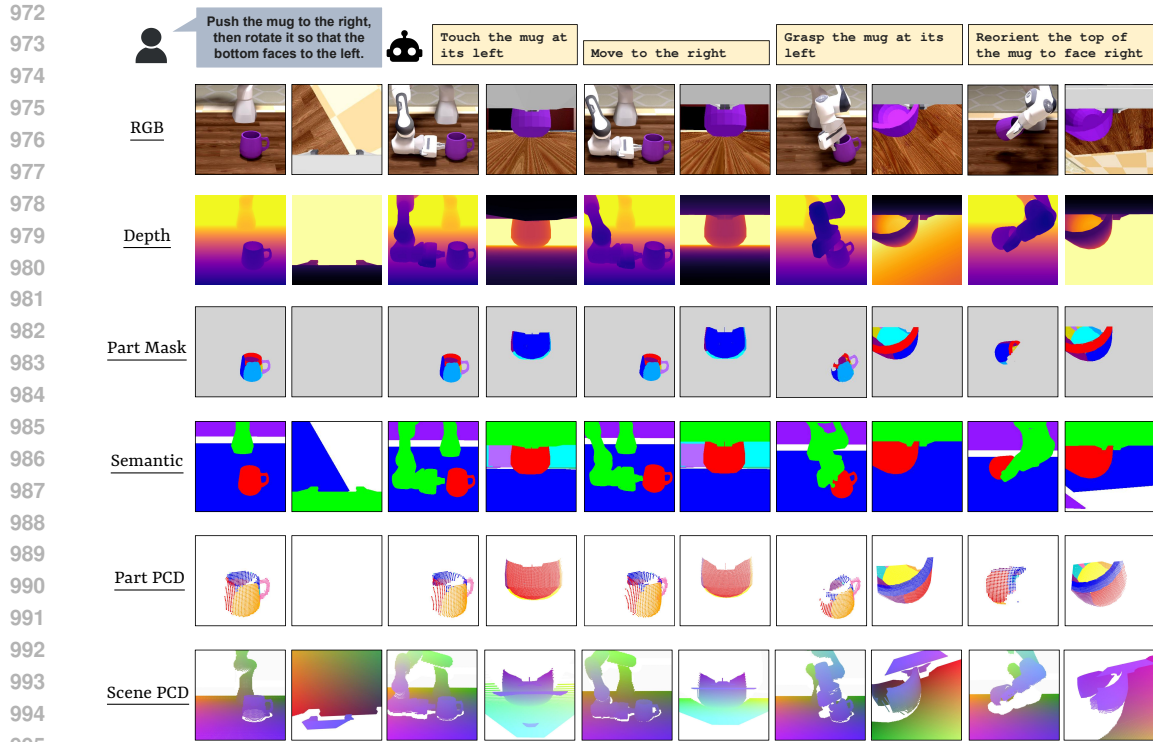


Figure 8: Different visual modalities in PartGym.

A.2.3 DETAILS OF PARTGYM

The aim of PartGym is to boost embodied AI research related to interaction with table-top object parts. PartGym support real-time rendering of different visual modalities (see Figure 8). In addition to the typical modalities like RGB, depth and object segmentation, PartGym also provides part-related visions like part masks and part point cloud, including spatial parts and semantic parts of an object. These part-related vision modalities are rendered by *PyBullet* (Coumans & Bai, 2016–2021) simulation engine using the ground-truth part assets given by *PartNet Mobility* (Xiang et al., 2020a; Mo et al., 2019; Chang et al., 2015).

Additionally, PartGym provides a framework to implement bi-level planning models for part-level manipulation tasks in simulation environments. It provides a template skill instruction generator, an oracle skill execution checker, as well as a systematical way to render part-related modalities shown in any skill instruction.

A.2.4 VISUALIZATION OF DIFFERENT TEST SPLITS

We provide the visualization of all 5 test sets in this section.

Move the bottle on the table away from me while holding top, then free it.



Figure 9: Left: Training set. Right: Test 1(OS).

Place gripper tip on the screw of the scissors.



Figure 10: Left: Training set. Right: Test 2(OI).

Grab the top of the mug and move it forwards.



Grab the handle of the mug and move it backwards.



Figure 11: Above: Training set. Below: Test 3(TP).

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

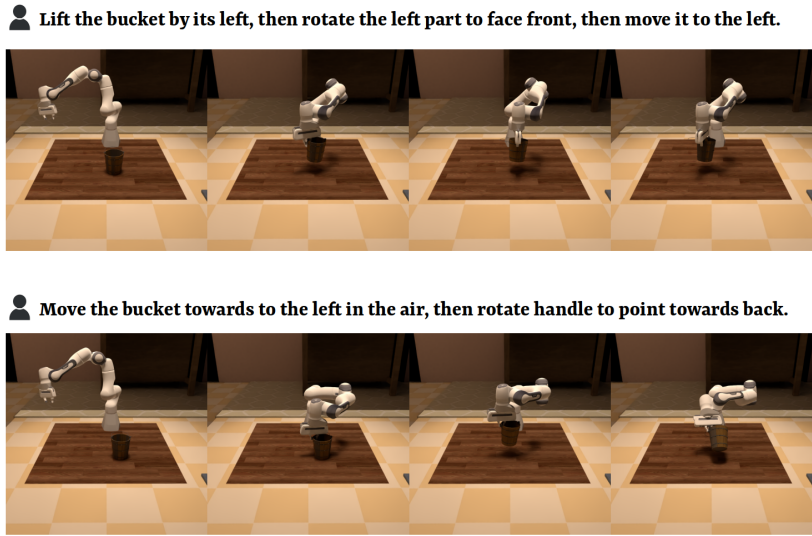


Figure 12: Above: Training set. Below: Test 4(TC).



Figure 13: Left: Training set. Right: Test 5(OC).

A.2.5 STATISTICS OF PARTINSTRUCT EPISODES

We provided detailed statistics about parts within each object type.

1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187

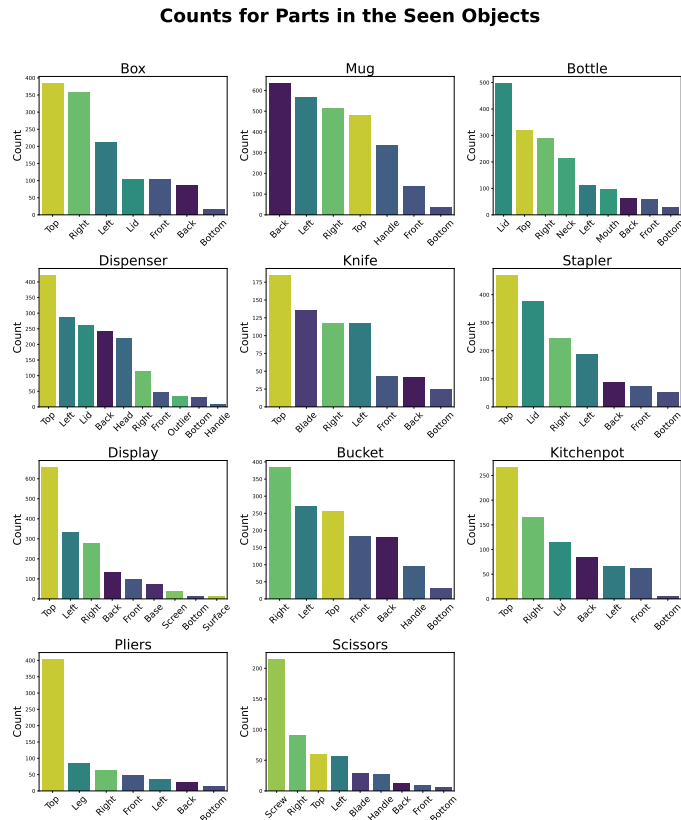


Figure 14: Parts in PartInstruct episodes, grouped by seen object types.

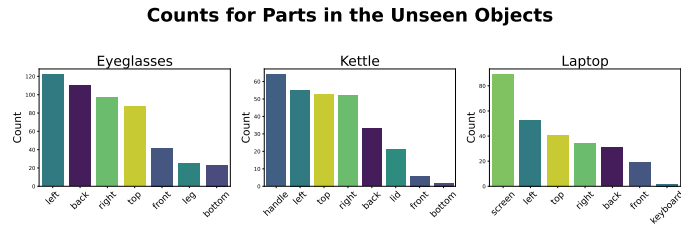


Figure 15: Parts in PartInstruct episodes, grouped by unseen object types.

A.3 SKILL AND OBJECT PART IMPACT STUDY

Here, we selected the rollout logs of the best-performing policy and analyzed the impact of different skill types and object parts. Specifically, we evaluated the success rate and failure causes for each skill and part.

The Success Rate was calculated by dividing the number of successful executions of each skill or part by the number of times it appeared in the skill chain. The Failure Cause was calculated by dividing the number of times a skill chain failed because of a specific skill or part by the total number of skill chain failures.

Table 14: Average success rate and failure cause for the three part-level skills

Skill	Grasp Object	Rotate Object	Touch Object
Success Rate (%)	53.55	18.18	54.55
Failure Cause (%)	43.51	6.11	16.79

Table 15: Average success rate and failure cause for selected parts.

Part	Blade	Left	Neck	Top	Screen	Mouth	Bottom
Success Rate (%)	46.67	45.10	66.67	52.78	60.00	66.67	30.00
Failure Cause (%)	4.60	13.79	1.15	24.14	2.30	0.00	0.00

Part	Handle	Leg	Lid	Front	Right	Back	Screw	Head
Success Rate (%)	64.29	33.33	63.64	29.03	36.49	41.03	0.00	16.67
Failure Cause (%)	4.60	2.30	4.60	5.75	21.84	9.20	1.15	3.45

A.4 IMPLEMENTATION DETAILS

A.4.1 TRAINING DETAILS IN END-TO-END POLICY LEARNING

We trained the baseline models, including Diffusion Policy (DP) (Chi et al., 2023), 3D Diffusion Policy (DP3) (Ze et al., 2024), and Act3D (Gervet et al., 2023), from scratch. For RVT2 (Goyal et al., 2023) and Octo (Team et al., 2024b), we implemented both fine-tuning of the pretrained models and training from scratch on our dataset. All trained models are using vision modalities from a static-view camera put with the same extrinsics in the workspace, as well as the real-time robot states information. Experiments were conducted on cluster nodes of A100 or H100 using Distributed Data Parallel (DDP). Training from scratch generally took about two days, while fine-tuning required one day.

DIFFUSION POLICY (DP)

We train CNN-based DP from scratch on our dataset. The action prediction horizon is set to 16 steps, with an observation horizon of 2 steps and action steps of 8. The input RGB images are cropped to a size of 76×76 . For language instructions, we use a pre-trained T5-small language encoder to obtain a language embedding of 512 dimensions. This language embedding is then concatenated with other features to form the final feature representation.

3D DIFFUSION POLICY (DP3)

The DP3 model is trained under a similar setup as DP, with an action prediction horizon of 16 steps, an observation horizon of 2 steps, and action steps of 8. For the point cloud observations, we use an input size of 1024 points, which are downsampled from the original point cloud using the Iterative Farthest Point Sampling algorithm (Qi et al., 2017). The language instructions are processed in DP3 following the same approach as in DP.

1242 ACT3D

1243

1244

1245 Act3D takes an image input size of 256×256 . The action prediction horizon is set to 6 steps, and the
 1246 observation horizon is 1 step. Following the raw work (Gervet et al., 2023), we use ResNet50(He
 1247 et al., 2016) as the vision encoder, and use CLIP (Radford et al., 2021) embeddings for vision-
 1248 language alignment. For 3D action map generation, the number of "ghost" points is set to be 10,000,
 1249 with a number of sampling level of 3.

1250

1251 OCTO

1252

1253 For fine-tuning, we use the released checkpoint of the `octo-base-1.5` model and fine-tune its
 1254 output head for 20,000 iterations. We use both the static view camera and the wrist view camera.
 1255 The input image sizes are 256×256 for the static view and 128×128 for the wrist view. The
 1256 window size is set to 2 steps, and the action horizon is set to 16 steps.

1257

1258

1259 RVT2

1260

1261 To adapt RVT2 in our benchmark settings, we first convert the depth map from the static camera
 1262 view into a point cloud in the camera coordinates, then apply camera extrinsic to transfer the point
 1263 cloud into the world coordinates, where the action heat maps will be generated and apply supervision
 1264 in. The action prediction horizon is chosen to be 6 steps, and the observation horizon is set to be 1
 1265 step.

1266

1267 A.4.2 RESULTS OF PRE-TRAINED GENERALIST POLICIES

1268

1269 We selected several popular generalist policy, including RT-1, Octo and OpenVLA, and evaluated
 1270 their zero-shot performance on our test sets. For RT-1, we followed the implementation of Open X-
 1271 Embodiment project and used the released `rt_1_x_tft_trained_for_002272480_step` check-
 1272 point for inference. For Octo, we used `octo-base-1.5` model, following the same setup as de-
 1273 scribed in A.4.1. For OpenVLA, we used the pretrained model `openvla-7b`. We followed the
 1274 same evaluation protocol as other baselines, and our results shows that these generalist policies fail
 1275 to achieve any success on our test sets.

1276

1277

1278 A.4.3 IMPLEMENTATION DETAILS IN BI-LEVEL PLANNING

1279

1280 IMPLEMENTATION OF THE HIGH-LEVEL TASK PLANNER

1281

1282 The high-level task planner features a skill inference mechanism that leverages comprehensive con-
 1283 textual information—including user task instructions, previously executed skill chains, and real-time
 1284 state data such as vision and pose information—to determine the next appropriate action. Recall
 1285 that the high-level task planner updates the skill instruction once every n steps. Here n is de-
 1286 termined by the average number of steps typically required for each skill in the training dataset.
 1287 Specifically, we use 130 for `grasp_obj`, 30 for `move_gripper`, 68 for `touch_obj`, 40 for
 1288 `release_obj`, and 22 for `rotate_obj`. Once the execution counter reaches these average val-
 1289 ues, the VLM is prompted to infer the subsequent skill based on the current state. This design
 1290 ensures that the decisions are grounded in both historical and real-time data. In addition, the plan-
 1291 ner incorporates an exception handling measure to maintain output consistency and reliability. Any
 1292 unacceptable terms generated by the VLM—such as directional indicators that out of our defini-
 1293 tion, will be normalized to their prescribed equivalents. Also, despite embedding all acceptable part
 1294 names within the prompt, a dedicated mechanism cross-references any inferred part names against
 1295 a stored list of valid part names. For VLM baselines, we use `gemini-1.5-flash-002` and
`gemini-2.0-flash-exp` for Gemini models Team et al. (2024a), and `GPT4o` for the OpenAI
 model Islam & Moushi (2024).

1296 TRAINING OF THE LOW-LEVEL ACTION POLICIES

1297

1298

1299

1300 We trained the low-level action policy using skill instructions retrieved from our training data, as-
 1301 suming the presence of an oracle planner that decomposes the overall task. Apart from the skill
 1302 instructions, the training setup remains identical to the end-to-end learning approach described in
 1303 A.4.1

1304

1305

1306

1307

1308

1309

1310 PART GROUNDING AND TRACKING

1311

1312

1313 We selected `sam2_hiera_small` as our mask generation and tracking model due to its
 1314 fastest tracking time among all configurations of SAM 2. For language grounding, we chose
 1315 `Florence-2-large` as our Vision-Language Model. To evaluate the performance, we used the
 1316 rollout logs of *DP-S SAM2*.

1317

1318 The performance was assessed using two key metrics: Grounding Success, Intersection over Union
 1319 (IoU). Grounding Success is calculated as the ratio of successfully grounded parts to the total number
 1320 of parts during a task. A grounding is considered successful if: 1) after language grounding, the
 1321 prompt points given by the VLM consist of one positive and one negative point (to prompt SAM
 1322 2), and 2) the IoU of the generated mask is greater than zero. If either of these conditions is not
 1323 met, the grounding is deemed a failure. The IoU measures the overlap between the predicted mask
 1324 generated by SAM 2 and the ground-truth mask retrieved from PartGym environment, it is defined
 1325 as the area of intersection divided by the area of union of the predicted and true regions. The results
 1326 across different test sets are summarized in Table 16.

1326

1327

1328

1329

1330

1331 Table 16: Performance of part grounding and tracking across different test sets.

1332

Metric	Test1	Test2	Test3	Test4	All
Grounding Success (%)	25.26	35.73	36.87	15.06	27.58
IoU	0.15	0.18	0.19	0.25	0.20

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347 VLM PROMPTS

1348

1349 We provide implementation details on high-level task planner in this section. Below is the VLM
 prompts.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

Prompt Example

You are an expert at planning manipulation tasks. You will be given one task instruction for each manipulation task. Each task instruction can be divided into a chain of skill instructions. Your job is to infer the next skill instruction (you only need to output one immediate next skill instruction each time, even if the entire task requires multiple skills) for the robot to execute, based on the following information and the attached image (current rgb frame) without outputting any intermediate inference and explanation:

- **Task Instruction:** {user_input}
- **Executed Skill Instructions:** {executed_skill_instructions}
- **Gripper State:** gripper.state (The gripper is open when the value is around 0.04 and it is closed when the value is less or around 0.018.)
- **First and Last TCP Pose:** This is the tcp pose of the first frame: first_tcp_pose and this is the tcp pose of the last frame: last_tcp_pose. (The tcp format is np.r.[self.rotation.as_quat(), self.translation] where the first four elements are the rotation quaternion and the last three are the translation element. You can use the tcp poses difference to infer if the move_gripper skill is complete (the Euclidean distance of two tcp positions should be at least 0.05). Please analyze the images and determine whether the task is completed between the first and last frames. Focus on the position of the object relative to the gripper. A small gap between the gripper and the object can be allowed.)

Here is more relevant information about the task.

Skill Descriptions:

1. **grasp_obj:**
 - *Description:* This skill grasps an object by a specific part.
 - *Parameters:*
 - part_grasp:** The exact part of the object to be grasped. Must match the user's input (e.g., 'blade', 'lid').
 - *Format:* Grasp the {obj_class} at its {part_grasp}
2. **move_gripper:**
 - *Description:* This skill moves the gripper in a specified direction while optionally keeping an object grasped.
 - *Parameters:*
 - dir_move:** Direction to move the gripper. Can only be 'top', 'bottom', 'left', 'right', 'front', or 'back'.
 - *Format:* Move {dir_str}, where 'dir_str' is mapped from 'dir_move' by: - 'front' → 'forwards' - 'back' → 'backwards' - 'top' → 'upwards' - 'bottom' → 'downwards' - 'left' → 'to the left' - 'right' → 'to the right'
3. **rotate_obj:**
 - *Description:* This skill rotates an object in a specific direction based on a given part.
 - *Parameters:*
 - dir_rotate:** Direction to rotate the object. Must be one of 'top', 'bottom', 'left', 'right', 'front', 'back'.
 - part_rotate:** The part of the object that should be rotated.
 - *Format:* Reorient the {part_rotate} of the {obj_class} to face {dir_str}, where 'dir_str' is mapped from 'dir_rotate'.
4. **touch_obj:**
 - *Description:* This skill touches a part of an object.
 - *Parameters:*
 - part_touch:** The part of the object to be touched.
 - *Format:* Touch the {obj_class} at its {part_touch}
5. **release_obj:**
 - *Description:* This skill releases an object from the gripper.
 - *Parameters:* None.
 - *Format:* Release

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

Prompt Example (Continued)

Part Names:

- **Scissors:** blade, handle, screw, left, right, top, bottom, front, back
- **Kitchen Pot:** base body, lid, left, right, top, bottom, front, back
- **Laptop:** base frame, screen, touchpad, keyboard, screen frame, left, right, top, bottom, front, back
- **Eyeglasses:** base body, leg, left, right, top, bottom, front, back
- **Bucket:** handle, base body, left, right, top, bottom, front, back
- **Display:** base support, surface, frame, screen, left, right, top, bottom, front, back
- **Pliers:** base body, leg, outlier, left, right, top, bottom, front, back
- **Bottle:** mouth, lid, body, neck, left, right, top, bottom, front, back
- **Knife:** base body, translation blade, rotation blade, left, right, top, bottom, front, back
- **Stapler:** base body, lid, body, left, right, top, bottom, front, back
- **Kettle:** handle, base body, lid, left, right, top, bottom, front, back
- **Mug:** handle, body, containing things, left, right, top, bottom, front, back
- **Box:** rotation lid, base body, left, right, top, bottom, front, back
- **Dispenser:** base body, pressing lid, head, handle, outlier, left, right, top, bottom, front, back

Task Splitting Example:

Break down the task: Split the task instruction into individual steps. Example: "Move the box in the air towards the right while keeping touch of right, then put it down." Steps: (1) Grasp the box at its right, (2) Move upwards, (3) Move to the right, (4) Move downwards. Return only the next skill instruction in the specified format.

Notes:

- Do not modify or assume alternate names for object parts.
- The task sequence should follow the user's input as strictly as possible.
- Do not replace object parts with similar or inferred names.
- Only bucket, mug, and scissors have a part called handle. Do not infer the handle part name for other objects.