

ARNOLD: A Benchmark for Language-Grounded Task Learning with Continuous States in Realistic Scenes

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Understanding continuous object states and task goals is essential for
2 task planning since they are generally not discrete in the real world. However,
3 most previous task learning benchmarks assume discrete (*e.g.*, binary) object states,
4 making it barely applicable to transfer the policy from the simulated environment
5 to the real world. Moreover, the trained robot’s ability to follow human
6 instructions based on grounding the actions and states is limited. To address such
7 challenges, we present ARNOLD, a benchmark that evaluates language-grounded
8 task learning with continuous states in realistic 3D scenes. ARNOLD consists of 8
9 language-conditioned manipulation tasks that require an in-depth understanding of
10 continuous object states and policy learning for continuous goals. To encourage
11 language-instructed learning, we provide template-generated demonstrations with
12 language descriptions. We will benchmark the task performances with state-of-
13 the-art language-conditioned policy learning algorithms. We will release ARNOLD
14 and host challenges to promote future research in embodied AI and robotics.

15 **Keywords:** Grounded task learning, Continuous states, Simulated environment

16 1 Introduction

17 One key capability emerged from the long evolution of human beings is the ability to ground language.
18 Such capability enables humans to describe and learn concepts, execute tasks, and also communicate
19 with others. With the recent culmination of grounding concepts in images [Radford et al., 2021;
20 Kamath et al., 2021; Saharia et al., 2022], few have studied the grounding of actions [Shridhar et al.,
21 2022; Zheng et al., 2022]. Considering how humans understand object status and relate language
22 instructions to the physical world, a natural question to ask is: How can we equip robotics system the
23 same capability to ground language instructions to task executions in a physical world?

24 There exists several challenges that cause non-trivial difficulties in the learning of robotic systems.
25 First, robot tasks are highly dependent on the detailed scene information, especially geometry
26 information, layouts, and visual appearances[Xing et al., 2021]. This requires perception modules
27 of robots to acquire detailed scene understanding, extracting useful geometrical features for task
28 execution. This challenge is further aggravated by the various combinations of different scene
29 configurations, including novel appearances, objects, and spatial positions. It is therefore critical for
30 robotic systems to learn and generalize skills of a task to novel scene configurations.

31 Moreover, an essential capability of human lies in the precise understanding of desired goal states.
32 Although humans often refer to goals with simple descriptions (*e.g.*, a cup *half* filled, a door *fully*
33 opened, *etc.*), what we understand is actually the status of physical properties (*e.g.*, half the volume,
34 pulled to 180°, *etc.*). Given this abstraction, it is exceedingly difficult for robots to learn the accurate
35 goal state from abstracted task descriptions, not to mention more abstract language descriptions
36 that refer to a implicit range of continuous object states (*e.g.*, *a bit* of coffee, *slightly* open, *etc.*).
37 This further requires the understanding of continuous object states. To take into account such subtle

38 capabilities, robot systems need to maintain a mapping from language instruction to an estimate of
39 desired goal states before task execution.

40 A first step towards the aforementioned robot learning problems is on building robot simulation
41 systems that facilitate learning and simulation. In fact, recent years have witnessed significant
42 progress in simulated environments that facilitate grounded task learning [Das et al., 2018; Shridhar
43 et al., 2020; Mees et al., 2021; Zheng et al., 2022]. However, these benchmarks exhibit several
44 limitations that prevent the robots from achieving capabilities that can function in real world: 1)
45 assuming discrete (*e.g.*, binary) object states and perfect motor control, therefore ignoring, low-level
46 geometry and dynamics of the object [Szot et al., 2021; Srivastava et al., 2022; Ehsani et al., 2021], not
47 requiring in-depth physical state understanding or fine-grained manipulation skills; 2) not grounding
48 instructions to precise states [Zheng et al., 2022; Shridhar et al., 2022], omitting the challenging
49 grounding problem from language to specific states in a continuous spectrum; 3) performing tasks
50 in simple and clean environments instead of in scenes spatially constrained by various surrounding
51 objects and visually disturbed by diverse textured backgrounds [Kumar and Todorov, 2015; Lin et al.,
52 2020; Zheng et al., 2022].

53 To better address the grounded task learning problem in a realistic setting, we introduce, ARNOLD,
54 a new benchmark for grounding task description languages to **continuous robot actions** and **con-**
55 **tinuous object states** in a **photo-realistic** and **physical-realistic** interactive environment. ARNOLD
56 is built on top of Nvidia Issac Sim, which provides accurate physics simulation and state-of-the-art
57 rendering. We use a hybrid human-template based approach to synthesize data for eight different
58 tasks. Each task features continuous robot motion with friction-based grasping and object state
59 manipulations. These tasks require different motor skills, including but not limited to grasping,
60 pushing, pulling, and pouring liquid into a container for 40 unique objects in 20 different scenes.
61 Then we pair each demonstration with a template based language instruction that describes task goals.

62 Different from prior work, ARNOLD tests the agent’s ability to generalize to unseen object states
63 through understanding language instructions and continuous object states.

64 In summary, our grounded continuous task learning benchmark ARNOLD has the following contribu-
65 tions:

66 • A set of eight different **language-conditioned manipulation** tasks featuring different motor skills
67 and diverse object states.

68 • A **photo-realistic** and **physical-realistic** 3D simulation environment with **continuous states** for
69 different objects and fluids.

70 2 Related Work

71 **Simulation Environment for Embodied AI.** There has been a large number of works that simulate
72 indoor household activities for training and evaluating AI agents [Kolve et al., 2017; Puig et al., 2018;
73 Li et al., 2021; Szot et al., 2021]. Most of these simulators only simulate preconditions and post
74 effects of agent actions or use simplified state representations. Instead, our work simulates continuous
75 states for articulated objects and simulates fluids at the particle level.

76 Some works use simplified abstract discrete action space. Abstract discrete action space can reduce
77 the task’s difficulty. However, agents trained in this setting are not aware of the low-level geometry
78 and dynamics of the objects, which would restrict their possibility of transferring to the real world.
79 For example, grasping is often simplified by attaching a nearby object to the gripper [Ehsani et al.,
80 2021; Shridhar et al., 2020; Srivastava et al., 2022; Szot et al., 2021], or through contact [Li et al.;
81 James et al., 2020]. By contrast, we control robots with 7-DOF continuous control and friction-based
82 grasping powered by state-of-the-art physics engines (PhysX5.0).

83 Among works that simulate continuous object state change, Gao et al. [2019] predict an action verb
84 giving task goals. However, they do not manipulate object states in a fine-grained manner and only
85 assume binary object states. Softgym[Lin et al., 2020] is an object manipulation benchmark that

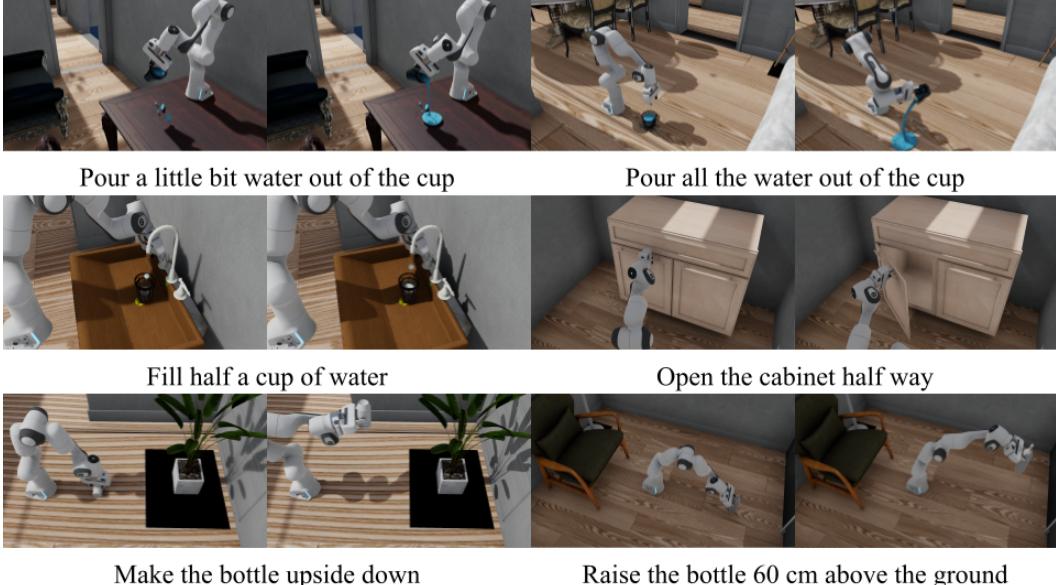


Figure 1: Examples of tasks in ARNOLD. For accomplishing these tasks, it requires visual recognition, text understanding, as well as diverse cognitive skills.

86 provide realistic simulation of deformable objects. However, there is a lack of diversity in terms of
 87 objects and scenes. In contrast, ARNOLD provide a wide variety of scenes and objects. In addition,
 88 as Isaac Sim is a part of Omniverse, authoring more complex assets can be easily managed through
 89 other applications within the omniverse universe.

90 **Language Conditioned Manipulation.** Relating human language to robot actions has been of
 91 interest in recent research [Mees et al., 2021; Lynch and Sermanet, 2020; Shridhar et al., 2020; Zheng
 92 et al., 2022]. However, the environments in these works either lack realistic physics or do not have
 93 realistic scenes where the agent’s surroundings will constrain its motion, and different scene objects
 94 might occlude the agent’s viewpoint. Most importantly, prior work aims to ground human languages
 95 to static object properties, such as colors and shapes. By contrast, ARNOLD provides instructions for
 96 continuous object states.

97 **Continuous state understanding.** Some recent research tries to predict object states [Liu et al.,
 98 2017; Nagarajan and Grauman, 2018]. However, the object states are discrete rather than continuous.
 99 More recently, [Weng et al., 2021] tried to predict object states in a continuous spectrum. However,
 100 they do not involve manipulating objects from an arbitrary starting state to the desired state, and it
 101 only involves articulated objects. In addition, they do not model the language grounding process.

102 **3 Arnold Benchmark**

103 The aim of ARNOLD is to evaluate the learning of language-conditioned continuous control policy
 104 over a diverse range of physical-realistic and photo-realistic 3D scenes . In this setting, an agent
 105 needs to have the following capabilities:

- 106 1. Understand the goal state of the current object from a human language instruction.
- 107 2. Estimate current object states from multiple camera inputs.
- 108 3. Propose and actuate motion plans based on physics over long-horizon.

109 Concretely, for an instruction like pouring half of water out of the cup, the agent needs to understand
 110 the instruction and know what the cup should look like when it’s half full from multi-model sensors
 111 and reasoning. Then it needs to further reason over the best sequence of actions to grasp the cup and

| Benchmark | Alfred | Maniskill | Calvin | Behavior | KitchenShift | VLmbench | Ours |
|--------------|--------|-----------|--------|----------|--------------|----------|------|
| Language | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ | ✓ |
| Multi Camera | ✗ | ✓ | ✓ | ✗ | ✗ | ✓ | ✓ |
| Fluid | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ |
| Physics | ✗ | ✓ | ✓ | ✗ | ✓ | ✓* | ✓ |
| Continuous | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Scenes | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✓ |
| Robot | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

Table 1: **Comparison with other benchmarks.** ARNOLD features continuous control with **robots** over continuous object states with a large number of human demonstrations in photo-realistic scenes. Each task in ARNOLD is specified by a human language instruction. ARNOLD also leverages advanced physics simulations powered by PHYSX5.0 to simulate articulated bodies and fluids. **Language:** task goal is specified by a human language instruction. **Multi-Camera:** robot is equiped with multiple cameras. **Fluid:** advanced fluid simulation. **Physics:** realistic physics simulation. * RLbench based benchmarks use a simplified grasping. **Continuous:** object state is continuous. **Scene:** task are performed with scene background. **Robot** : perform actions with real robots.

112 pour the water out. Finally it needs to know when and how to stop to meet the goal object state it
 113 determined previously.

114 3.1 Simulation Environment

115 3.1.1 Environment Components

116 **Simulation Platform.** ARNOLD is based on NVIDIA Isaac Sim[Makoviychuk et al., 2021], a robotics
 117 simulation application that enables the creation of photo-realistic and physically-accurate virtual
 118 environments to develop robots. The physics simulation in ARNOLD is based on PhysX 5.0 and the
 119 photo-realistic rendering is powered by GPU-enabled ray tracing. See Figure 1 for some examples of
 120 simulation and rendering. As demonsrated by Li et al., the rendering effect provided by Isaac Sim is
 121 much more realisitc than any existing simulators.

122 **Scenes and Robots.** The scenes in ARNOLD are based on Fu et al. [2021], a large-scale synthetic
 123 indoor scenes dataset. As a result, ARNOLD scenes have professionally designed layouts and are
 124 populated by high-quality 3D models. Each scene contains a 7-DOF Franka Emika Panda manipulator
 125 with a parallel gripper.

126 **Fluids.** The fluids in ARNOLD are simulated using GPU-accelerated position-based-dynamics(PBD)
 127 method Macklin and Müller [2013] through omniverse. Then we perform surface construction
 128 through marching cubes Lorensen and Cline [1987] to obtain the final fluid rendering effect. For
 129 faster rendering, we make the surface construction step optional.

130 **Articulated Objects.** In addition to objects provided by Isaac Sim, we also get object parts from
 131 open-sourced dataset Kolve et al. [2017]; Xiang et al. [2020]. We perform modifications to a few
 132 object meshes, e.g. changing materials and adding covering meshes to cabinets and drawers, for
 133 more natural appearance. We also perform convex decomposition to create realistic collision for each
 134 object. Moreover, to ensure physically-realistic simulation, we assign physics parameters to objects,
 135 including weight and friction for rigid-body objects, and cohesion, surface tension and viscosity for
 136 fluids. We chose these parameters based on Mu et al. [2021] and human operator feedback.

137 3.1.2 Observation and Action Space of the Robot

138 **Action space.** The agent can control the 7 joints and the gripper of the manipulator. In addition, the
 139 agent can control robot end-effector through the provided motion planner. Since the focus of this
 140 work is object manipulation, we do not allow the robot to navigate the room.

141 **Observation space.** The robot has five different cameras views around it. One is on top of the robot,
 142 one is on the opposite side of the robot, one is on the left of the robot and two are on the robot gripper.
 143 This setting is to avoid view occlusion problem during manipulation. Each camera provides RGB-D

144 input. By default, each camera has a resolution of 256×256 . Users can render arbitrary resolution
 145 through the replay. In addition, we provide robot joint positions and velocities. To support learning,
 146 we provide additional state information available to access if needed: robot base position, object
 147 positions, object rotations, object part semantic mask.

148 **3.2 Task Design**

149 **3.2.1 Task Definition**

150 In ARNOLD, each task $t \in T$ is defined by a tuple $(t_o \in O, t_a \in A_o, t_b \in B_s, t_g \in G_s)$. Here O is
 151 the set of object types the robot needs to manipulate (e.g. cabinet, drawer, glass). A_o is the set of
 152 attributes we care for object o (e.g. joint angle of the cabinet, water level in the container). B_s is the
 153 set of initial attribute values of the object being manipulated, e.g. the initial joint angle of the cabinet
 154 or the amount of water in a cup at the beginning. G_s is the set of required goal values of the object.

| Task Type | Goal variations | success condition |
|-----------------|----------------------------|-------------------|
| Pick up object | 10, 20, 30, 40 cm above | ± 5 cm |
| Reorient object | 45, 135, 180 degrees | ± 15 degree |
| Open cabinet | 25, 50, 75, 100 % open | ± 10 % |
| Open drawer | 25, 50, 75, 100 % open | ± 10 % |
| Close cabinet | 25, 50, 75, 100 % open | ± 10 % |
| Close drawer | 25, 50, 75, 100 % open | ± 10 % |
| Pour water | 25, 50, 75, 100 % of water | ± 10 % |
| TransferWater | 20, 40, 60, 80 % of water | ± 10 % |

Table 2: Task Description for eight different tasks. Each task features four different goal state variations specified by human language. For each task, the success signal will be triggered if the object state is within the success tolerance for two seconds. The tolerance threshold is developed based on a pilot study where we ask humans to use a gamepad to finish the task.

155 **3.2.2 Task Types**

156 In ARNOLD, there are nine types of tasks. Here we define each task type.

157 **PickUpObject.** $T_{pick} = \{t | t_o = bottle \wedge t_a = height \wedge (t_b < t_g)\}$. This task requires the agent
 158 to find the object, pick it up and raise it to the desired height. This task is meaningful in real-life
 159 human-robot collaboration where the human instructs the robot to raise the object to a certain height
 160 so he can receive the object from the robot.

161 **ReorientObject.** $T_{reorient} = \{t | t_o = bottle \wedge t_a = orientation\}$. To execute this task, the agent
 162 needs to find the object, pick it up, and reorient it to the desired angle. A real-life example of this task
 163 can be getting ketchup out of the glass bottle, which requires the agent to reorient the bottle to the
 164 desired angle.

165 **OpenCabinet/Drawer.** $T_{open} = \{t | t_o = \{Cabinet, Drawer\} \wedge t_a = jointVal \wedge (t_b < t_g)\}$.
 166 For this task, the agent needs to estimate the current joint position of the cabinet/drawer and take
 167 appropriate actions based on human language. For cabinets, the motions are constrained by a revolute
 168 joint. For drawers, the motions are constrained by a prismatic joint.

169 **CloseCabinet/Drawer.** $T_{open} = \{t | t_o = \{Cabinet, Drawer\} \wedge t_a = jointVal \wedge (t_b > t_g)\}$. This
 170 task is the reverse of the previous task. It is relatively easy compared to other tasks; however, it
 171 requires the agent to anticipate the consequences of its actions correctly. If the agent exerts a large
 172 force to push the object initially, the drawer will be closed completely.

173 **PourWater.** $T_{pour} = \{t | t_o = \{Container\} \wedge t_a = waterLevel \wedge (t_b > t_g)\}$. To successfully
 174 execute this task, the agent needs to pick up the cup and pour a desired amount of water out of the cup.
 175 Therefore, the agent needs to estimate the amount of water in the cup and take appropriate actions.

176 **TransferWater.** $T_{transfer} = \{t | t_o = \{\text{Container}\} \wedge t_a = \text{waterLevel} \wedge (t_b < t_g)\}$. The agent
 177 needs to pick up the cup and transfer a specified amount of water from one cup to another for this
 178 task. This task is the most difficult task of all. First, the agent needs to know where to pour water
 179 without knocking over the container cup. Then the agent needs to estimate how much water is in the
 180 cup and stop at the right time. To avoid ambiguity, two cups are the same.

181 Detailed task descriptions are in table 2

182 3.2.3 Evaluation Metrics

183 **Task Success.** Each task is parametrized by a tuple as indicated in section 3.2.1. Task success is
 184 defined as in table 2. The goal state is satisfied when the object stays in a tolerance threshold, as
 185 defined below:

$$|t_a(i) - t_g| < \epsilon_t, \quad (1)$$

186 where ϵ_t is the success tolerance for each task, as shown in Table 2.

187 If the object states satisfies the goal state for 2 seconds, then we declare the task as success. For
 188 example, consider the task: "pour half a cup of water out of the cup". The agent succeeds if, there are
 189 40% – 60% of water particles in the cup for two seconds within the episode limit.

190 3.2.4 Task split

191 To allow generalization on novel scenes and objects, we split the dataset into four folds: 1) training,
 192 2) novel scene textures, where objects are seen during training but scene textures are different, 3)
 193 novel objects, where scenes textures are seen during training but not objects, and 4) novel language
 194 states.

195 3.3 Dataset Collection

196 3.3.1 Mission Definition

197 We call any particular robot task instance a mission. Each mission is defined by a tuple $(t \in T, i \in$
 198 $I_o, e \in E)$, which specifies the setting combinaiton of task T , object I_o , and environment E . So in
 199 a mission, the robot needs to perform a task t in a specific environment e using an object instance
 200 i . Moreover, we still allow the position of the robot and object in the scene to vary in a mission to
 201 create more diversity.



Figure 2: Multi View Camera Example

202 3.3.2 Human Demonstrations and Instructions

203 To understand how humans complete the task, we collect 2400 human demonstrations with 7 operators
 204 on all task types.

205 We created 400 missions (t, i, e) consisting of various combinations of object instances, goal state
 206 values, and scenes. For each mission, the annotator would further specify two relative positions
 207 between the robot and the object. For each position setting, the annotator would record two human
 208 demonstrations. So for each mission, we will record three human demonstrations.

209 We enable operators to teleoperate with the robot via the Xbox controller during human recording. To
 210 avoid view occlusion with other objects and furniture in the scene, we also allow operators to change
 211 the camera viewpoint using the controller joystick.

212 In addition, for each human demonstration collected, we sample 3 template-based language commands
213 based on the language generation engine to describe the target object state. Note that the initial states
214 are not specified in the command, thus requiring the agent to understand the current object state from
215 observations.

216 Human demonstrations are served as a sanity check to ensure that all tasks can be completed. However,
217 human demonstrations are very noisy. It poses too many difficulties for current models to extract
218 meaningful information. To this end, we proposed a keypoint-based task template inspired by [Zheng
219 et al., 2022; Shridhar et al., 2022]

220 3.3.3 Keypoint Based Task Templates

221 We manually mark keypoints around objects as in [Zheng et al., 2022] to generate noise-free tra-
222 jectories.. Each task composes of several stages. At each stage, a motion planner drives the robot
223 toward a keypoint. When the stage condition is satisfied, task planner will move on to the next stage.
224 We designed a set of keypoint-based task planners for all objects and tasks in our benchmark. To
225 generate feasible locations in the scene for objects and robots, we reuse human-annotated (object,
226 robot) location pairs.

227 4 Baseline Model

228 Recently, there are two proposed baselines for language conditioned robot manipulation tasks.
229 6D CLIPort [Zheng et al., 2022] assumes a known task template and uses a top-down projection of
230 fused point cloud to predict keypoint on the images space, then translates the projection back to the
231 world space. To handle the missing height and roll from the top-down projection, they directly use a
232 regressor to predict those two quantities. PerAct [Shridhar et al., 2022], instead, directly uses a voxel
233 representaton to predict keypoints for different stages of the action.

234 The above two methods assume a strong spatial prior, making the prediction data-efficient and robust..

235 5 Limitations

236 Regarding the simulation environment, even though we have put realism of the simulation as one
237 of our key focuses, there is still some gap to the real world due to the precision limit of physical
238 simulation and rendering. Rudin et al. [2022] has shown that sim-to-real is possible with policies
239 trained under Isaac Sim. We plan to perform sim-to-real experiment in future work.

240 Some of our task settings are not very natural as we would see in the real world, though, due to
241 limited 3D assets and limited manpower. For example, in real-life, people generally pour water into
242 the sink instead of on the ground or the table except for cleaning purposes. But such a setting should
243 not limit the robots learning the critical manipulation tasks, which is the focus of this work.

244 Besides, robots are limited in their type variation. We collect data and perform the experiment on
245 only one type of robot. There is no variation in robot arm length and degree of freedom. We thought
246 changing robot types might pose too much difficulty for our tasks.

247 Regarding the language instruction design, we did not use natural language generated by human
248 annotators in this work. Template-based language is not as diverse as natural language. Therefore, we
249 cannot achieve a good coverage for human language descriptions to object states. At the same time,
250 as indicated by Jansen [2020], language instructions generated by annotators, even though diverse,
251 contain significant errors. So our instructions will also lack some real-world noise.

252 Besides, we allow operators to perturb the robot’s relative positions against the object to avoid
253 the overfitting indicated in Mees et al. [2021]. Operator-generated positions only cover a limited
254 percentage of robot possible locations. More general settings should enable the robot to navigate.
255 But this work focuses on object manipulation instead of navigation, so we put aside this complexity.
256 We will consider adding it to our future work.

257 **6 Conclusion and Future Work**

258 We presented ARNOLD, a new benchmark for language-conditioned continuous control over physical-
259 realistic and photo-realistic scenes. ARNOLD contains eight different tasks requiring diverse skills.
260 We further supplies human demonstrations as well as template generated demonstrations. ARNOLD
261 tries to bridge the gap between object state understanding and object manipulation for modern robotics
262 systems.

263 **References**

- 264 A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin,
265 J. Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.
- 267 A. Kamath, M. Singh, Y. LeCun, G. Synnaeve, I. Misra, and N. Carion. Mdetr-modulated detection for
268 end-to-end multi-modal understanding. In *Proceedings of the IEEE/CVF International Conference
on Computer Vision*, pages 1780–1790, 2021.
- 270 C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. Denton, S. K. S. Ghasemipour, B. K. Ayan, S. S.
271 Mahdavi, R. G. Lopes, et al. Photorealistic text-to-image diffusion models with deep language
272 understanding. *arXiv preprint arXiv:2205.11487*, 2022.
- 273 M. Shridhar, L. Manuelli, and D. Fox. Cliport: What and where pathways for robotic manipulation.
274 In *Conference on Robot Learning*, pages 894–906. PMLR, 2022.
- 275 K. Zheng, X. Chen, O. C. Jenkins, and X. E. Wang. Vlmbench: A compositional benchmark for
276 vision-and-language manipulation. In *Proceedings of the Neural Information Processing Systems
Track on Datasets and Benchmarks*, 2022.
- 278 E. Xing, A. Gupta, S. Powers, and V. Dean. Kitchenshift: Evaluating zero-shot generalization of
279 imitation-based policy learning under domain shifts. In *NeurIPS 2021 Workshop on Distribution
Shifts: Connecting Methods and Applications*, 2021.
- 281 A. Das, S. Datta, G. Gkioxari, S. Lee, D. Parikh, and D. Batra. Embodied question answering. In
282 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–10,
283 2018.
- 284 M. Shridhar, J. Thomason, D. Gordon, Y. Bisk, W. Han, R. Mottaghi, L. Zettlemoyer, and D. Fox.
285 Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In *Proceedings of
286 the IEEE/CVF conference on computer vision and pattern recognition*, pages 10740–10749, 2020.
- 287 O. Mees, L. Hermann, E. Rosete-Beas, and W. Burgard. Calvin: A benchmark for language-
288 conditioned policy learning for long-horizon robot manipulation tasks. *arXiv preprint
289 arXiv:2112.03227*, 2021.
- 290 A. Szot, A. Clegg, E. Undersander, E. Wijmans, Y. Zhao, J. Turner, N. Maestre, M. Mukadam, D. S.
291 Chaplot, O. Maksymets, et al. Habitat 2.0: Training home assistants to rearrange their habitat.
292 *Advances in Neural Information Processing Systems*, 34, 2021.
- 293 S. Srivastava, C. Li, M. Lingelbach, R. Martín-Martín, F. Xia, K. E. Vainio, Z. Lian, C. Gokmen,
294 S. Buch, K. Liu, et al. Behavior: Benchmark for everyday household activities in virtual, interactive,
295 and ecological environments. In *Conference on Robot Learning*, pages 477–490. PMLR, 2022.
- 296 K. Ehsani, W. Han, A. Herrasti, E. VanderBilt, L. Weihs, E. Kolve, A. Kembhavi, and R. Mottaghi.
297 Manipulathor: A framework for visual object manipulation. In *Proceedings of the IEEE/CVF
298 conference on computer vision and pattern recognition*, pages 4497–4506, 2021.

- 299 V. Kumar and E. Todorov. Mujoco haptix: A virtual reality system for hand manipulation. In *2015*
300 *IEEE-RAS 15th International Conference on Humanoid Robots (Humanoids)*, pages 657–663.
301 IEEE, 2015.
- 302 X. Lin, Y. Wang, J. Olkin, and D. Held. Softgym: Benchmarking deep reinforcement learning for
303 deformable object manipulation. *arXiv preprint arXiv:2011.07215*, 2020.
- 304 E. Kolve, R. Mottaghi, W. Han, E. VanderBilt, L. Weihs, A. Herrasti, D. Gordon, Y. Zhu, A. Gupta, and
305 A. Farhadi. Ai2-thor: An interactive 3d environment for visual ai. *arXiv preprint arXiv:1712.05474*,
306 2017.
- 307 X. Puig, K. Ra, M. Boben, J. Li, T. Wang, S. Fidler, and A. Torralba. Virtualhome: Simulating
308 household activities via programs. In *Proceedings of the IEEE Conference on Computer Vision*
309 and Pattern Recognition, pages 8494–8502, 2018.
- 310 C. Li, F. Xia, R. Martín-Martín, M. Lingelbach, S. Srivastava, B. Shen, K. Vainio, C. Gokmen,
311 G. Dharam, T. Jain, et al. Igibson 2.0: Object-centric simulation for robot learning of everyday
312 household tasks. *arXiv preprint arXiv:2108.03272*, 2021.
- 313 C. Li, C. Gokmen, G. Levine, R. Martín-Martín, S. Srivastava, C. Wang, J. Wong, R. Zhang,
314 M. Lingelbach, J. Sun, et al. Behavior-1k: A benchmark for embodied ai with 1,000 everyday
315 activities and realistic simulation. In *6th Annual Conference on Robot Learning*.
- 316 S. James, Z. Ma, D. R. Arrojo, and A. J. Davison. Rlbench: The robot learning benchmark & learning
317 environment. *IEEE Robotics and Automation Letters*, 5(2):3019–3026, 2020.
- 318 X. Gao, R. Gong, T. Shu, X. Xie, S. Wang, and S.-C. Zhu. Vrkitchen: an interactive 3d virtual
319 environment for task-oriented learning. *arXiv preprint arXiv:1903.05757*, 2019.
- 320 C. Lynch and P. Sermanet. Language conditioned imitation learning over unstructured data. *arXiv*
321 *preprint arXiv:2005.07648*, 2020.
- 322 Y. Liu, P. Wei, and S.-C. Zhu. Jointly recognizing object fluents and tasks in egocentric videos. In
323 *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- 324 T. Nagarajan and K. Grauman. Attributes as operators: factorizing unseen attribute-object composi-
325 tions. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 169–185,
326 2018.
- 327 Y. Weng, H. Wang, Q. Zhou, Y. Qin, Y. Duan, Q. Fan, B. Chen, H. Su, and L. J. Guibas. Capra:
328 Category-level pose tracking for rigid and articulated objects from point clouds. In *Proceedings of*
329 *the IEEE/CVF International Conference on Computer Vision*, pages 13209–13218, 2021.
- 330 V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin,
331 A. Allshire, A. Handa, et al. Isaac gym: High performance gpu-based physics simulation for robot
332 learning. *arXiv preprint arXiv:2108.10470*, 2021.
- 333 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:
334 3d furnished rooms with layouts and semantics. In *Proceedings of the IEEE/CVF International*
335 *Conference on Computer Vision*, pages 10933–10942, 2021.
- 336 M. Macklin and M. Müller. Position based fluids. *ACM Transactions on Graphics (TOG)*, 32(4):
337 1–12, 2013.
- 338 W. E. Lorensen and H. E. Cline. Marching cubes: A high resolution 3d surface construction algorithm.
339 *ACM siggraph computer graphics*, 21(4):163–169, 1987.
- 340 F. Xiang, Y. Qin, K. Mo, Y. Xia, H. Zhu, F. Liu, M. Liu, H. Jiang, Y. Yuan, H. Wang, et al. Sapien: A
341 simulated part-based interactive environment. In *Proceedings of the IEEE/CVF Conference on*
342 *Computer Vision and Pattern Recognition*, pages 11097–11107, 2020.

- 343 T. Mu, Z. Ling, F. Xiang, D. Yang, X. Li, S. Tao, Z. Huang, Z. Jia, and H. Su. Maniskill: Generalizable
344 manipulation skill benchmark with large-scale demonstrations. *arXiv preprint arXiv:2107.14483*,
345 2021.
- 346 M. Shridhar, L. Manuelli, and D. Fox. Perceiver-actor: A multi-task transformer for robotic manipu-
347 lation. *arXiv preprint arXiv:2209.05451*, 2022.
- 348 N. Rudin, D. Hoeller, P. Reist, and M. Hutter. Learning to walk in minutes using massively parallel
349 deep reinforcement learning. In *Conference on Robot Learning*, pages 91–100. PMLR, 2022.
- 350 P. A. Jansen. Visually-grounded planning without vision: Language models infer detailed plans from
351 high-level instructions. *arXiv preprint arXiv:2009.14259*, 2020.

352 **A Appendix for ARNOLD: A Benchmark for Language-Grounded Task**
 353 **Learning with Continuous States in Realistic Scenes**

354 **A.1 Parsing Assets into Universal Scene Description (USD) format**

355 Parsing assets into Omniverse is parsing graphic data into USD files. The goal is to make it easy for
 356 developers to read their assets (such as scene files, articulation bodies, animations e.t.c.) and have
 357 that information available when needed. This means that users can use a USD file with a bunch of
 358 different assets and then tell OMNIVERSE which ones are important for your game at any given time
 359 without having to go through each one individually.

360 **A.1.1 Working with 3D-Front dataset**

361 The 3D-Front dataset is originally a collection of synthetic indoor scenes highlighted with professional
 362 design and a large number of rooms with high-quality textured 3D models. The source dataset contains
 363 three main parts: models, scene files, and textures. The dataset consists of more than tens of thousands
 364 of room layouts with thousands of furnished objects.

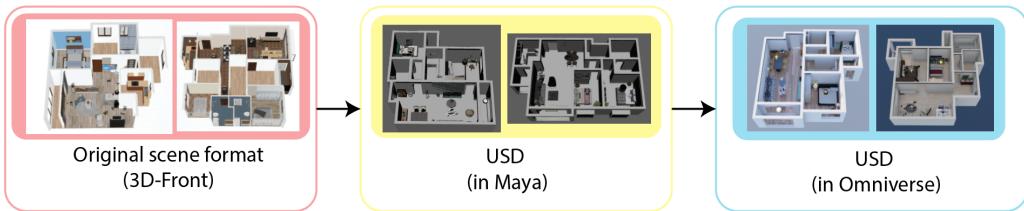


Figure 3: Pipeline for parsing scenes. First, we pre-process the original 3D-Front scene data. Then, we build an automatic pipeline to load the scene layout into Autodesk Maya with custom designs. Finally, we convert the layout file into USD format and deploy it in Omniverse.

365 To parse the 3D-Front dataset into usd format, we apply the following steps:

- 366 • Parsing the original scene files (.json) into a data frame containing the mesh and furniture
 367 information;
- 368 • Using Maya MEL script to load scenes into Autodesk Maya;
- 369 • Applying Maya and Omniverse converter to save the scenes into USD format.

370 **A.1.2 Working with Articulation bodies**

371 The articulation bodies in our application mainly come from the SAPIEN. To parse the original
 372 articulation bodies (.urdf) into USD format, we apply OMNIVERSE ISAAC SIM. The build-in tool
 373 within it allows to transfer urdf files into desired format.

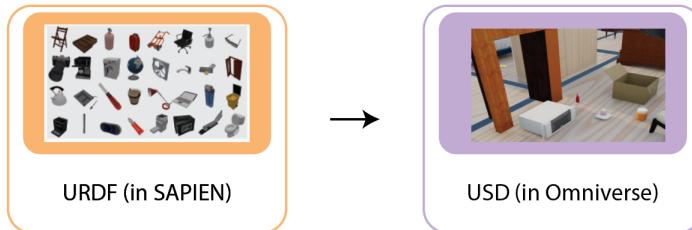


Figure 4: Pipeline for parsing articulation bodies. To convert the original articulation bodies in URDF format, we refine the built-in functionalities in Omniverse Isaac-Sim and partly modify the assets manually. Then we save the parsed objects in to USD format.

374 **A.2 Working with robotics**

375 Robotics is a science that deals with the design, construction and operation of robots. The field
376 encompasses many different disciplines such as computer science, electrical engineering, mechanical
377 engineering and others. In Omniverse, robotics is built-in with sufficient details. Users may make it
378 stronger by customizing the designs or even making their own robots from scratch.



Figure 5: Robots in Omniverse Isaac Sim. The Omniverse platform provides a lot of features to simulating virtual robotics. It helps researchers with the tools they need to build robust, physical-realistic simulations.

379 **A.3 Data Collection**

380 **A.3.1 Robot Teleoperation**

381 **Frame Transformation.** For data collection, we enable the operator to control the robot and the
382 camera using an Xbox controller. We assume that the controller input is given in the robot base
383 frame. As displayed in Figure 6, the base frame is define as the frame where the X axis is originated
384 from the robot's base and pointing to the object and the Y axis is the up axis. The controller input is
385 transformed to the world frame at each time step.

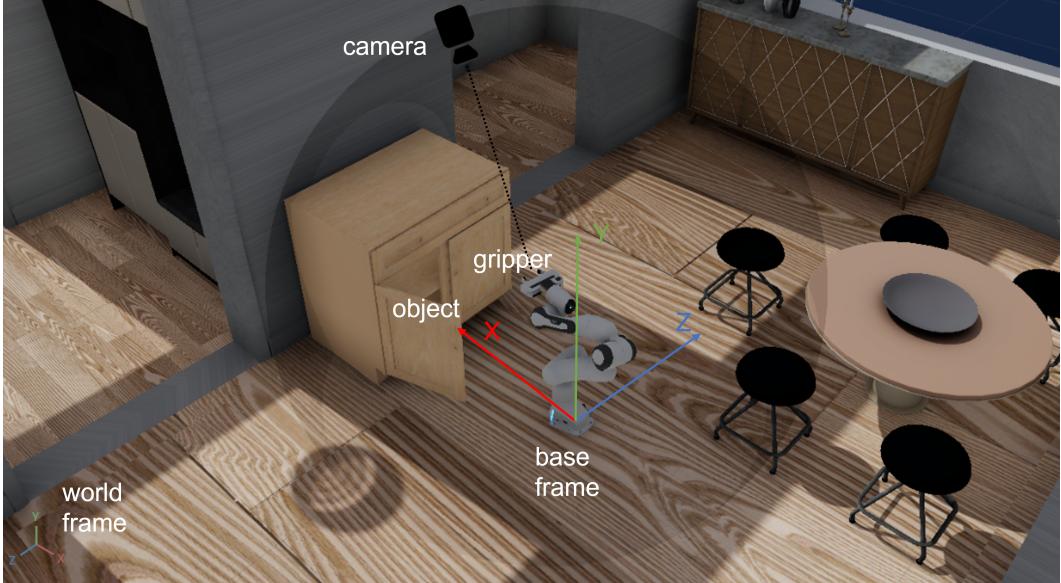


Figure 6: Frames and camera

386 **Teleoperation.** For robot teleoperation, the operator can change the position and rotation of the
387 end effector and toggle the gripper (Figure 7). In addition to controlling the height and horizontal
388 movement of the end effector, operators can change its rotation via two control modes (i.e. joint
389 position control and end effector rotation control). The joint position control mode is intuitive since it

390 allows the operator to directly change the positions of joints A1 (shoulder joint), A6 (forearm joint)
 391 and A7 (wrist joint). The end effector rotation control allows operators to rotate the end effector
 392 around the X, Y, Z axis of the base frame while maintaining its position, and thus is a more general
 393 way of changing the rotation. The operator can switch between these two modes during the data
 394 collection.

395 **Camera control.** We enable operators to move their viewing camera freely in the virtual environment
 396 to avoid occlusions in cluttered scenes so that they can continuously monitor the end effector and
 397 the object. This is done by a spherical camera control (Figure 6), where the camera is continuously
 398 facing the end effector and the operator can move it around a sphere centered at the end effector. The
 399 radius of the sphere can also be adjusted for a more clear view.



Figure 7: Controller Mapping



Figure 8: User Interface.

| Task Name | Template | Example |
|----------------|-------------------------------|--|
| PickObject | <v><o><s> | raise the bottle 20cm above the ground |
| ReorientObject | <v><o><s> | rotate the bottle 90 degrees from the up axis |
| OpenCabinet | <v><position><o><s> | pull the left cabinet half closed |
| CloseCabinet | <v><position><o><s> | close the top cabinet completely |
| OpenDrawer | <v><position><o><s> | open the bottom drawer midway |
| CloseDrawer | <v><position><o><s> | push the bottom drawer entirely closed |
| PourWater | <v><s> water <preposition><o> | pour all the water out of the cup |
| TransferWater | <v><s> water <preposition><o> | transfer almost half of the water to the glass |

Table 3: language generation template. We sample each attribute based during the generation process based on the proposed template. For each trajectory, we sample 3 different language descriptions. (Where, <v>, <o>, and <s> stand for *verb*, *object*, and *state* respectively.)

400 A.3.2 Procedure

401 The human data collection process is done using a user interface (UI) Figure 8, which is implemented
 402 as an extension of Isaac Sim. For each mission (t, i, e) , the annotator is given a uniquely sampled
 403 tuple specifying the object instance, goal state values and scene. During the process, the annotator
 404 also needs to specify the transformation of the robot and object relative to the scene.
 405 The annotator starts by moving the object group containing the robot and the object instance to
 406 an appropriate anchor place in the scene and clicking on the "Record house + anchor" button for
 407 recording. The annotator can then adjust the relative transformation between the object and the
 408 robot and click on the "Record Robot" button. For each mission, the annotator is required to change
 409 the transformation twice. After clicking on the "load mission and record" button, the annotator
 410 can control the robot to manipulate the object. In each trial, the goal is to change the attribute
 411 t_a of object t_o instance i from t_b to t_g . Two trials are needed for each mission with the adjusted
 412 transformation. Once a mission's success condition is met, the UI displays a "Task Success" message
 413 and the annotator can click the "Stop" button to stop the recording. Finally, the annotator can replay
 414 the recorded trajectory by clicking the "Load mission & Replay" button. To proceed to the next
 415 mission, the annotator can click on the "Next" button.

416 A.3.3 Instructions Generation

417 For each trajectory collected, we sample 3 template-based language command based on the language
 418 generation engine to describe the object state as described in Table 3. For example, "pour 50% of
 419 water out of the cup." Note that the initial states are not specified in the command, thus requiring the
 420 agent to understand the current object state from observations.

421 A.3.4 Dataset Statistics

422 In ARNOLD, We sample actions from a Xbox robot controller at 120 Hz resulting in 2488 trajectories
 423 or around 4,400,000 frames of image and action pairs. The recorded trajectory length is about 10.2
 424 hours, with a median trajectory length of 13 seconds, an average of 14.8 seconds, a maximum of 91.6
 425 seconds, and a minimum of 3.2 seconds.

426 Then for each (human, object) pair, we generate an additonal trajectory through templates. Note, not
 427 all scenaiors that can are slovable by human are slovable by motion planning.

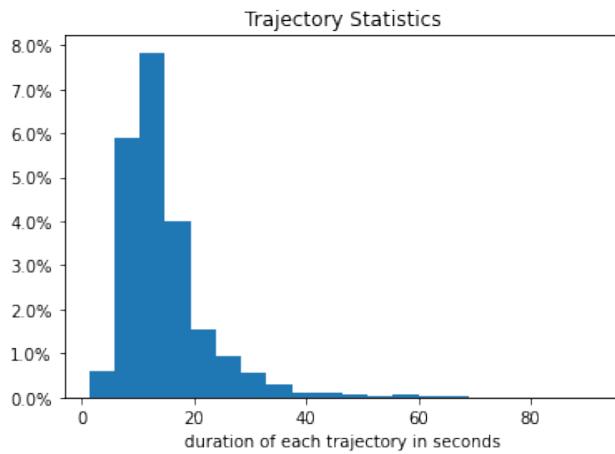


Figure 9: Trajectory Length Distributions

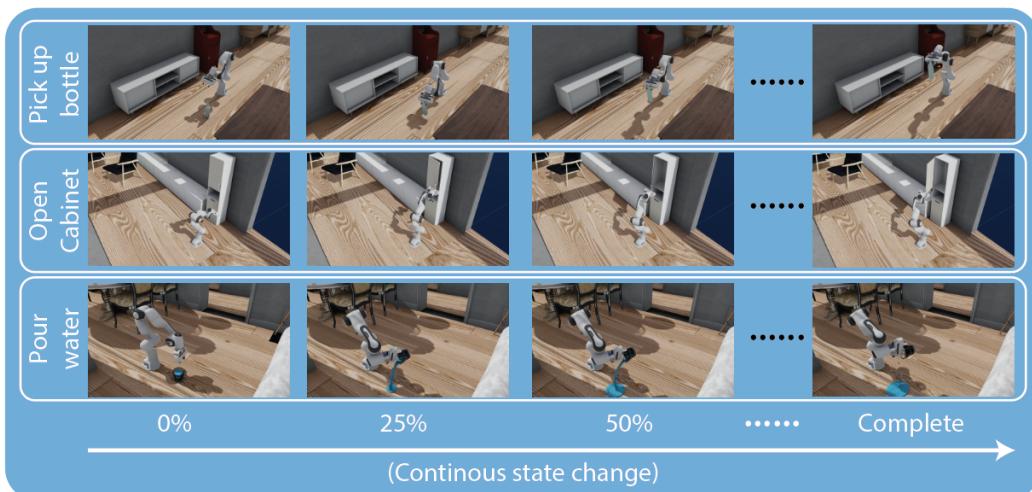


Figure 10: Continous control and goal state

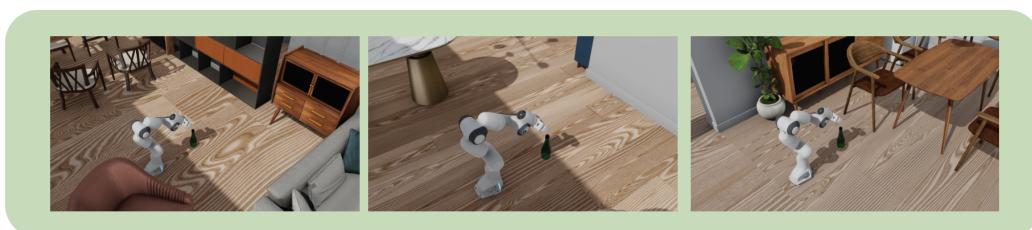


Figure 11: Scene variations

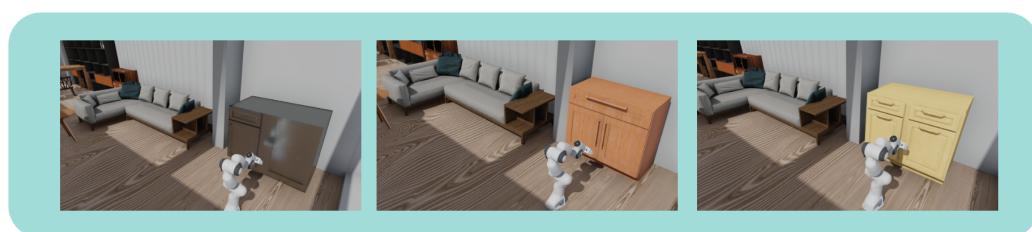


Figure 12: Object variations

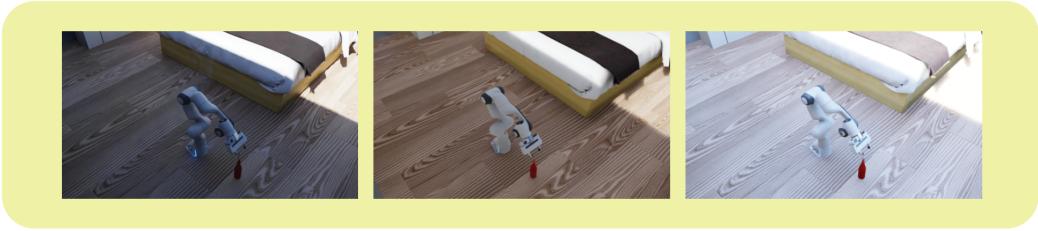


Figure 13: Lighting variations

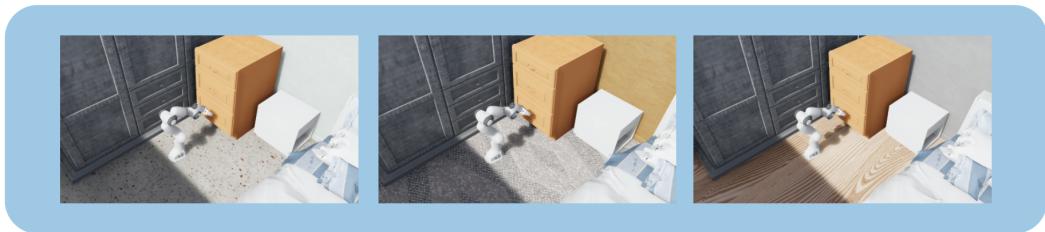


Figure 14: Material variations



Figure 15: Object variations