

# ROBOVERSE: A Unified Benchmark for Scalable and Generalizable Vision-Language Robotic Manipulation

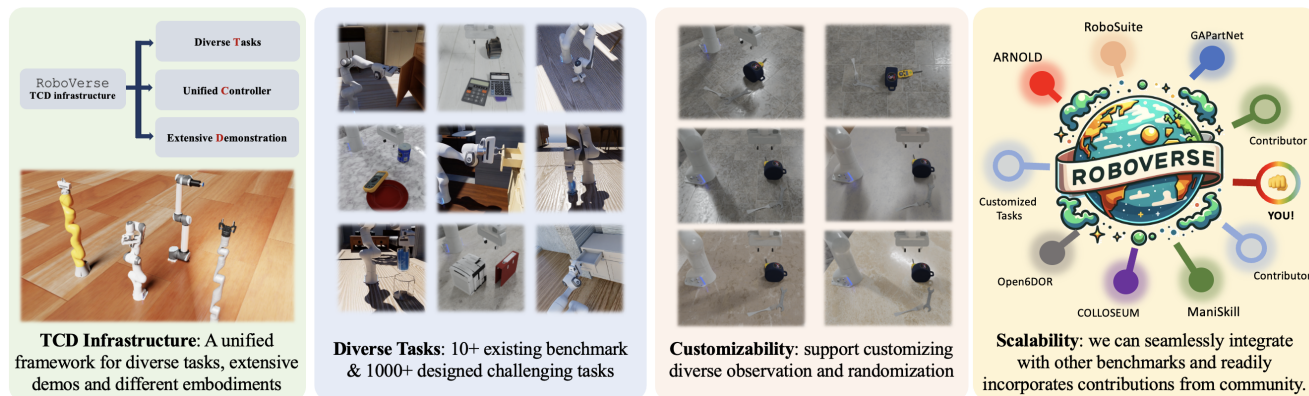


Figure 1. **Overview of RoboVerse benchmark.** We provide a unified infrastructure for robotic manipulation in simulation environments. This design unifies diverse tasks, extensive demonstrations, and different robot embodiments in existing robotic manipulation benchmarks. We further enrich the available robotic manipulation demonstration by scaling existing tasks with domain randomization and incorporating newly designed tasks. Our benchmark and dataset exhibit remarkable flexibility for each task, allowing for utilization across different observation modalities, diverse randomization strategies, and scalability with the joint efforts of the robotics community via an easy-to-use coding pipeline.

## Abstract

The importance of diverse, high-quality datasets is underscored by their role in training foundational models, especially in fields like natural language processing and computer vision. However, scaling up data and models for robotics presents unique challenges due to the confinement of prior models to specific datasets and domains and the limitations inherent in collecting diverse real-world demonstrations. To overcome these limitations, we propose leveraging simulators as an alternative. Simulators can generate vast, diverse datasets and allow for flexible manipulation of various elements, such as observation representations and action formats, thereby offering a scalable and adaptable approach for training robotic models. To this end, we propose RoboVerse benchmark, in which we provide a unified infrastructure for diverse tasks, extensive demonstrations, and different robot embodiments. We also collect a large-scale dataset merging both existing benchmarks and newly designed tasks. Furthermore, our framework exhibits remarkable flexibility, allowing for utilization across different observation modalities, diverse randomization strategies, and scalable data augmentation.

## 1. Introduction

Recent advancements in foundational models highlight the growing importance of comprehensive and high-quality datasets in improving model performance and generalization capabilities. However, directly adapting such data scaling effects to robotics research faces several significant challenges in data collection. Predominantly, with data from different sources utilizing different input modalities (e.g., RGB images, point clouds, etc.) and robot embodiments (e.g., Franka Emika, UR10e, etc.), setting a universal standard for both data representation and task is difficult. Consequently, linking research findings across different experimental settings for a cohesive conclusion is challenging.

To address this limitation, prior works have attempted to collect large-scale manipulation demonstrations in both real-world and simulation environments. In real-world settings, the RT series [1–4, 14] have RGB recordings of robot manipulation at the cost of extensive data collection efforts. Simulation-based benchmarks such as ManiSkill [10, 12] and RoboSuite [16] collect demonstrations of specific tasks in different simulation environments, making it challenging to transfer policies learned in different benchmarks.

Recognizing the substantial effort needed to collect real-world data and the lack of unification in robot simulation benchmarks, our RoboVerse benchmark encompasses the following appealing features:

- 048 • **Unified Structure:** We devise a unified task structure, 099  
049 dataset format, and evaluation system to support a wide 100  
050 range of tasks, complemented by a suite of tools to facili- 101  
051 tate task design and demonstration generation. 102
- 052 • **Flexibility and Diversity:** The flexibility of our bench- 103  
053 mark allows for effortless customization of new tasks, 104  
054 observation representations, and action representations to 105  
055 suit specific needs. We unify existing demonstrations in 106  
056 simulation environments and also collect a large amount 107  
057 of tasks and demonstrations with rich annotations for 108  
058 downstream policy learning. 109
- 059 • **Comprehensive Modalities:** Utilizing simulation envi- 110  
060 ronments, we provide both multi-view RGB recordings, 111  
061 and point clouds as well as detailed language annotations 112  
062 catering to a variety of tasks. 113
- 063 • **Scalability:** We provide extensive APIs to make our 114  
064 benchmark scalable for robotics manipulation research, 115  
065 enabling the seamless addition of new tasks, demonstra- 116  
066 tions, and the training of new models.

067 Utilizing RoboVerse, we can thoroughly evaluate 117  
068 the performance and generalization capabilities of existing 118  
069 methods across input modalities, tasks, and robot embodi- 119  
070 ments through both interpolation and extrapolation. Addi- 120  
071 tionally, with our integrated vision-language-action demon- 121  
072 strations, we can craft a versatile robotic manipulation pol- 122  
073 icy for diverse tasks and complex scenarios. 123

## 074 2. RoboVerse 124

075 We propose RoboVerse, a comprehensive and multi-task 125  
076 benchmark, developed using the IsaacSim simulator [13]. 126

### 077 2.1. Benchmark 127

078 **Unified Infrastructure.** We introduce our Task-Controller- 131  
079 Demonstration (TCD) Infrastructure for the RoboVerse 132  
080 benchmark. We have developed a base class for tasks, 133  
081 which includes all necessary functions and variables spe- 134  
082 cific to each task. Upon defining a task, we inherit from 135  
083 this base class, customizing configurations and environ- 136  
084 ment setup functions accordingly. Additionally, we have 137  
085 designed and implemented the controller infrastructure to 138  
086 manage embodiment, serving as an intermediary between 139  
087 the demonstration and the environment. It's worth not- 140  
088 ing that our controller is adaptable, accommodating vari- 141  
089 ous embodiments and standardizing their control mecha- 142  
090 nisms. Furthermore, we have collected extensive demon- 143  
091 strations for each task within our infrastructure, providing 144  
092 ample support for their utilization, re-rendering, replay, and 145  
093 evaluation. We have also established a standardized format 146  
094 for storing demonstrations across all tasks. 147

095 **Task.** Our RoboVerse benchmark integrates three key 148  
096 components: (1) pick and place, (2) articulated object ma- 149  
097 nipulation, and (3) complex manipulation tasks. Specifi-  
098 cally, we incorporate tasks, benchmarks, and demonstra-

tions from various sources into our benchmark, includ-  
ing: (1) ManiSkill [10, 12], (2) SceneDiffuser [11], (3)  
GAPartNet [8], (4) PartManip [7], (5) ARNOLD [9], (6)  
Open6DOR [5], (7) COLOSSEUM [15], (8) SAGE[6]. Fur-  
thermore, we customize certain specific tasks to address  
shortcomings identified in previous benchmarks. For ex-  
ample, we combine object pick and place with articulated  
object manipulation to create tasks such as “open the top  
drawer of the cabinet on the left of the table, and place the  
apple into it.” Additionally, we employ heuristics and re-  
inforcement learning algorithms to execute these tasks and  
gather demonstrations.

**Multiple Embodiment Support.** We support multiple  
embodiments in our RoboVerse benchmark, including  
Franka Emika FR3, Kinova Gen3, KUKA IIWA, Kinova  
Jaco, UR10e, and so on.

**Language Description.** We provide precise and detailed  
language descriptions for each task.

### 2.2. Demonstration 117

We gather large-scale demonstrations for the benchmark,  
each comprising task configurations, trajectories, language  
annotations, and other useful information. We introduce  
a unified representation for these demonstrations, focus-  
ing on the trajectory of the end-effector pose and gripper  
state. This standardized format enables the seamless reuse  
of demonstrations across various embodiments and scenes.  
Additionally, we offer comprehensive APIs for domain ran-  
domization, facilitating the creation of a more diverse and  
realistic dataset. Specifically, we currently support random-  
ization for (1) object colors, (2) ground plane, (3) lighting,  
(4) scene, (5) camera pose, (6) physical parameters, and (7)  
cross-embodiment.

For existing benchmarks with demonstrations, we di-  
rectly adopt theirs and transfer them to our simulator en-  
vironment, adjusting them to fit into the shared format. For  
other benchmarks or our newly designed tasks, we either  
utilize their existing policies, amalgamate several existing  
methods, or design heuristics to collect demonstrations.

## 3. Future Work 137

We are continuously expanding our benchmark and gather-  
ing more demonstrations to enhance our dataset. Leverag-  
ing these extensive demonstrations, we aim to quantitatively  
evaluate several key aspects crucial to the research commu-  
nity. Our goals include identifying optimal visual repre-  
sentations that ensure high performance and generalization,  
addressing data balancing challenges across various tasks,  
and assessing the model's generalization capabilities within  
and beyond its training distribution. Additionally, we plan  
to explore the contribution of simulation data to real-world  
applications, focusing on strategies for data balancing and  
identifying effective training paradigms.

## References

- 150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161  
162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205
- [1] Michael Ahn, Debidatta Dwibedi, Chelsea Finn, Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Karol Hausman, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, et al. Autort: Embodied foundation models for large scale orchestration of robotic agents. *arXiv preprint arXiv:2401.12963*, 2024. 1
- [2] Suneel Belkhole, Tianli Ding, Ted Xiao, Pierre Sermanet, Quon Vuong, Jonathan Tompson, Yevgen Chebotar, Debidatta Dwibedi, and Dorsa Sadigh. Rt-h: Action hierarchies using language. *arXiv preprint arXiv:2403.01823*, 2024.
- [3] 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.
- [4] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023. 1
- [5] Yufei Ding, Haoran Geng, Chaoyi Xu, Xiaomeng Fang, Jiazhao Zhang, Songlin Wei, Qiyu Dai, Zhizheng Zhang, and He Wang. Open6dor: Benchmarking open-instruction 6-dof object rearrangement and a vlm-based approach. In *First Vision and Language for Autonomous Driving and Robotics Workshop*. 2
- [6] Haoran Geng, Songlin Wei, Congyue Deng, Bokui Shen, He Wang, and Leonidas Guibas. Sage: Bridging semantic and actionable parts for generalizable manipulation of articulated objects. In *ICLR 2024 Workshop on Large Language Model (LLM) Agents*. 2
- [7] Haoran Geng, Ziming Li, Yiran Geng, Jiayi Chen, Hao Dong, and He Wang. Partmanip: Learning cross-category generalizable part manipulation policy from point cloud observations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2978–2988, 2023. 2
- [8] Haoran Geng, Helin Xu, Chengyang Zhao, Chao Xu, Li Yi, Siyuan Huang, and He Wang. Gapartnet: Cross-category domain-generalizable object perception and manipulation via generalizable and actionable parts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7081–7091, 2023. 2
- [9] Ran Gong, Jiangyong Huang, Yizhou Zhao, Haoran Geng, Xiaofeng Gao, Qingyang Wu, Wensi Ai, Ziheng Zhou, Demetri Terzopoulos, Song-Chun Zhu, et al. Arnold: A benchmark for language-grounded task learning with continuous states in realistic 3d scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 2
- [10] Jiayuan Gu, Fanbo Xiang, Xuanlin Li, Zhan Ling, Xiqiang Liu, Tongzhou Mu, Yihe Tang, Stone Tao, Xinyue Wei, Yunchao Yao, et al. Maniskill2: A unified benchmark for generalizable manipulation skills. *arXiv preprint arXiv:2302.04659*, 2023. 1, 2
- [11] Siyuan Huang, Zan Wang, Puhao Li, Baoxiong Jia, Tengyu Liu, Yixin Zhu, Wei Liang, and Song-Chun Zhu. Diffusion-based generation, optimization, and planning in 3d scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16750–16761, 2023. 206  
207  
208  
209  
210  
211
- [12] Tongzhou Mu, Zhan Ling, Fanbo Xiang, Derek Yang, Xuanlin Li, Stone Tao, Zhiao Huang, Zhiwei Jia, and Hao Su. Maniskill: Generalizable manipulation skill benchmark with large-scale demonstrations. *arXiv preprint arXiv:2107.14483*, 2021. 1, 2 212  
213  
214  
215  
216
- [13] NVIDIA. Isaacsim simulator. 2 217
- [14] Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anikait Singh, Anthony Brohan, et al. Open x-embodiment: Robotic learning datasets and rt-x models. *arXiv preprint arXiv:2310.08864*, 2023. 1 218  
219  
220  
221  
222
- [15] Wilbert Pumacay, Ishika Singh, Jiafei Duan, Ranjay Krishna, Jesse Thomason, and Dieter Fox. The colosseum: A benchmark for evaluating generalization for robotic manipulation. *arXiv preprint arXiv:2402.08191*, 2024. 2 223  
224  
225  
226
- [16] Yuke Zhu, Josiah Wong, Ajay Mandlekar, Roberto Martín-Martín, Abhishek Joshi, Soroush Nasiriany, and Yifeng Zhu. robosuite: A modular simulation framework and benchmark for robot learning. In *arXiv preprint arXiv:2009.12293*, 2020. 1 227  
228  
229  
230  
231