

UniDoorManip: Learning Universal Door Manipulation Policy Over Large-scale and Diverse Door Manipulation Environments

Yu Li^{*1} Xiaojie Zhang^{*1} Ruihai Wu^{*2}
 Zilong Zhang¹ Yiran Geng² Hao Dong^{†2} Zhaofeng He^{†1}
¹Beijing University of Posts and Telecommunications ²School of CS, Peking University

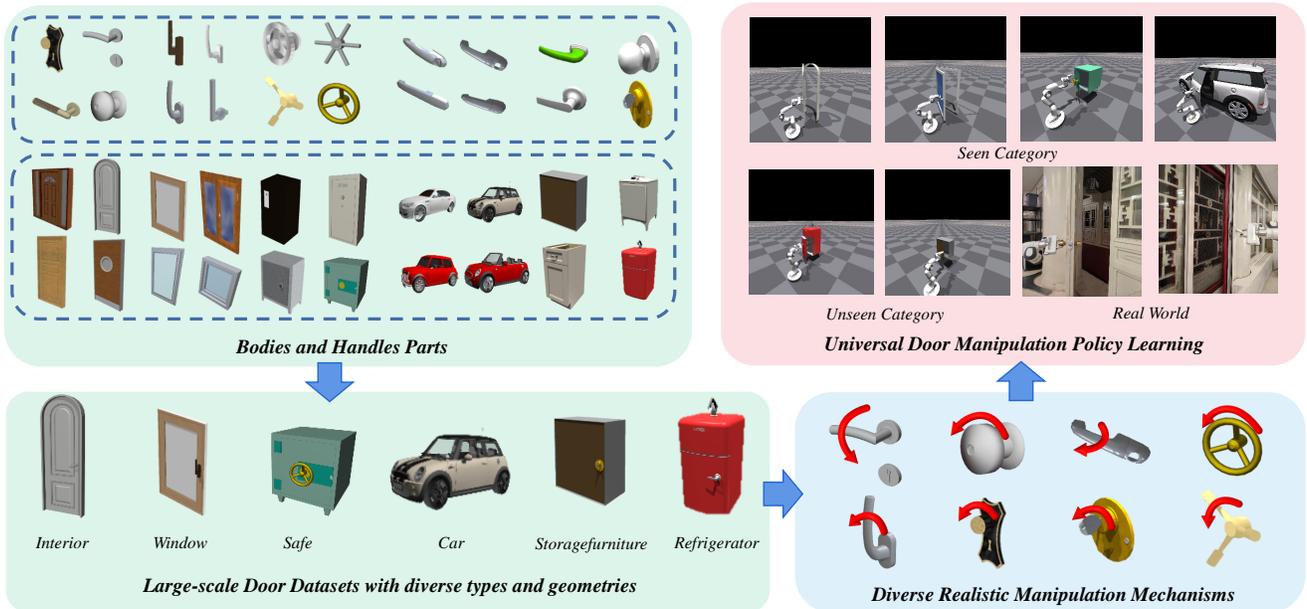


Figure 1. Our Proposed Environment, Dataset and Universal Manipulation Policy.

1. INTRODUCTION

Door manipulation holds significant importance due to the frequent need to open or close doors in various scenarios. While previous works have focused primarily on interior doors [10, 11], we aim to extend doors to a more general setting, *e.g.*, doors in windows, cars, safes, as illustrated in Figure 1. In the above broad scenarios, the door manipulation task covers doors with diverse types, geometries and manipulation mechanisms, which poses a great challenge to learn a universal door manipulation policy.

Due to the limited datasets and unrealistic simulation environments, previous works[1–3, 8, 13, 16] fail to achieve good performance across various doors. In this work, we **build a novel door manipulation environment** reflecting

different realistic door manipulation mechanisms, and further equip this environment with a **large-scale door dataset** covering 6 door categories with hundreds of door bodies and handles, making up thousands of different door instances as shown in Figure 1. Additionally, to better emulate real-world scenarios, we introduce a mobile robot as the agent and use the partial and occluded point cloud as the observation, which are not considered in previous works while possessing significance for real-world implementations. We conduct detailed comparisons between our proposed environment and dataset and others in Table 1, 2.

To learn a universal policy over diverse doors, we **propose a novel framework disentangling** the whole manipulation process into three stages, and integrating them by training in the reversed order of inference. Extensive experiments validate the effectiveness of our designs and demonstrate our framework’s strong performance. Code, data and

^{*}Equal contribution.

[†]Corresponding author.

Datasets	Int.			Win.			Car.			Saf.			Sto.			Ref.		
	B	H	CO	B	H	CO	B	H	CO	B	H	CO	B	H	CO	B	H	CO
AKB-48 [7]	-	9	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
PartNet-Mobility [15]	26	22	26	3	1	3	-	-	-	30	14	30	155	-	-	4	-	-
GAPartNet [3]	14	11	14	-	-	-	-	-	-	29	1	29	133	-	-	4	-	-
DoorGym [10]	-	20	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
Ours	57	96	5472	18	37	666	22	15	330	61	39	2379	160	8	1280	10	9	90

Table 1. **Statistic Comparisons Between Previous Dataset and Ours.** For category, **Int.**, **Win.**, **Car.**, **Saf.**, **Sto.**, **Ref.** respectively denote doors from Interior, Window, Car, Safe, StorageFurniture, Refrigerator. For asset number, **B**, **H**, **CO** indicate numbers of body, handle and composited object assets with the two parts.

Env.	Data.	Mob.	Latch.	Part.	Occ.
GAPartNet [3]	P + A				
W2A [8, 12, 13]	P			✓	
RLAfford [4]	P	✓			
PartManip [2]	G	✓		✓	
DoorGym [10]	D		✓	✓	
EnvAfford [14]	P			✓	✓
Ours	Ours	✓	✓	✓	✓

Table 2. **Comparison between Our Environment and Others.** For simplicity, **Data.**, **Mob.**, **Latch.**, **Part.** and **Occ.** respectively denote Dataset, Mobile Robot Arm, Latching Mechanism, Partial Observation and Occlusion in Observation. Besides, **P**, **A**, **G** and **D** respectively denote PartNet-Mobility, AKB-48, GAPartNet and DoorGym in Table 1.

videos are available on <https://unidoormanip.github.io/>.

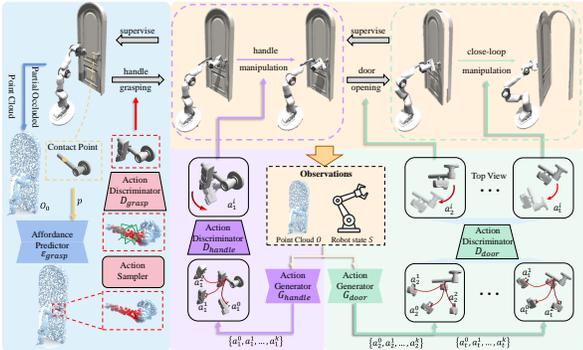


Figure 2. **Our Pipeline For The Framework.**

2. METHOD

As illustrated in Figure 2, we propose a novel framework that disentangles door manipulation into three distinct but related stages, each with a corresponding universal manipulation policy. We leverage conditioned training to train these policies, as they have inter-dependencies, and thus they can be integrated into a unified universal policy. In the first stage, we employ generalizable point-level visual affordance [5, 6, 9, 17] to propose stable grasp poses. In the second stage, we train a universal policy covering multiple handle manipulation mechanisms in our proposed realistic

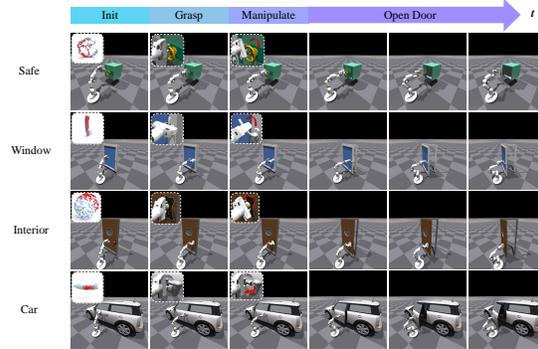


Figure 3. **Qualitative Results of Manipulation Sequence.**

Task	Pull Door					
	Train				Test	
Method	🚪	🚪	🚗	🚪	🚪	🚪
GAPartNet [3]+GT	0.62	0.88	0.41	0.44	0.52	0.26
DoorGym [10]	0.56	0.72	0.61	0.41	0.19	0.23
PartManip [2]	0.47	0.61	0.54	0.34	0.42	0.19
VAT-MART [13]	0.59	0.62	0.57	0.43	0.51	0.25
Ours w/o disentangle.	0.44	0.88	0.20	0.19	0.05	0.22
Ours w/o condition.	0.77	0.31	0.58	0.51	0.54	0.33
Ours w/o state.	0.73	0.59	0.16	0.36	0.45	0.37
Ours w/o mobile.	0.87	0.60	0.00	0.43	0.50	0.81
Ours	0.99	0.91	0.81	0.72	0.75	0.89

Table 3. **Quantitative Results of the Baselines and Ablations.**

environment. In the third stage, we train a policy to open doors with unlocked handles.

3. EXPERIMENTS

We conduct our experiments on the representative door manipulation tasks: **pull door**. The robot arm needs to pull the door until the door joint angle θ_d is larger than a threshold $thre_{door}$. Here, we set $thre_{door}$ to be 45° .

Figure 3 shows the whole manipulation sequence of our universal manipulation. We also compare our method with baselines and conduct an ablation study as shown in Table 3. Qualitative and quantitative results demonstrate that our universal policy can generalize over diverse categories, geometries and manipulation mechanisms.

References

- [1] Ben Eisner, Harry Zhang, and David Held. Flowbot3d: Learning 3d articulation flow to manipulate articulated objects. *arXiv preprint arXiv:2205.04382*, 2022. 1
- [2] 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
- [3] 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. 1, 2
- [4] Yiran Geng, Boshi An, Haoran Geng, Yuanpei Chen, Yaodong Yang, and Hao Dong. End-to-end affordance learning for robotic manipulation. *arXiv preprint arXiv:2209.12941*, 2022. 2
- [5] James J Gibson. The theory of affordances. *Hilldale, USA*, 1(2):67–82, 1977. 2
- [6] Suhan Ling, Yian Wang, Shiguang Wu, Yuzheng Zhuang, Tianyi Xu, Yu Li, Chang Liu, and Hao Dong. Articulated object manipulation with coarse-to-fine affordance for mitigating the effect of point cloud noise. *ICRA*, 2024. 2
- [7] Liu Liu, Wenqiang Xu, Haoyuan Fu, Sucheng Qian, Qiaojun Yu, Yang Han, and Cewu Lu. Akb-48: A real-world articulated object knowledge base. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14809–14818, 2022. 2
- [8] Kaichun Mo, Leonidas J. Guibas, Mustafa Mukadam, Abhinav Gupta, and Shubham Tulsiani. Where2act: From pixels to actions for articulated 3d objects. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6813–6823, 2021. 1, 2
- [9] Chuanruo Ning, Ruihai Wu, Haoran Lu, Kaichun Mo, and Hao Dong. Where2explore: Few-shot affordance learning for unseen novel categories of articulated objects. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. 2
- [10] Yusuke Urakami, Alec Hodgkinson, Casey Carlin, Randall Leu, Luca Rigazio, and Pieter Abbeel. Doorgym: A scalable door opening environment and baseline agent. *arXiv preprint arXiv:1908.01887*, 2019. 1, 2
- [11] Jiayu Wang, Shize Lin, Chuxiong Hu, Yu Zhu, and Limin Zhu. Learning semantic keypoint representations for door opening manipulation. *IEEE Robotics and Automation Letters*, 5(4):6980–6987, 2020. 1
- [12] Yian Wang, Ruihai Wu, Kaichun Mo, Jiaqi Ke, Qingnan Fan, Leonidas Guibas, and Hao Dong. Adaafford: Learning to adapt manipulation affordance for 3d articulated objects via few-shot interactions. *European conference on computer vision (ECCV 2022)*, 2022. 2
- [13] Ruihai Wu, Yan Zhao, Kaichun Mo, Zizheng Guo, Yian Wang, Tianhao Wu, Qingnan Fan, Xuelin Chen, Leonidas Guibas, and Hao Dong. VAT-mart: Learning visual action trajectory proposals for manipulating 3d articulated objects. In *International Conference on Learning Representations*, 2022. 1, 2
- [14] Ruihai Wu, Kai Cheng, Yan Zhao, Chuanruo Ning, Guanqi Zhan, and Hao Dong. Learning environment-aware affordance for 3d articulated object manipulation under occlusions. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 2
- [15] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao Jiang, Yifu Yuan, He Wang, et al. Sapien: A simulated part-based interactive environment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11097–11107, 2020. 2
- [16] Zhenjia Xu, He Zhanpeng, and Shuran Song. Umpnet: Universal manipulation policy network for articulated objects. *IEEE Robotics and Automation Letters*, 2022. 1
- [17] Yan Zhao, Ruihai Wu, Zhehuan Chen, Yourong Zhang, Qingnan Fan, Kaichun Mo, and Hao Dong. Dualafford: Learning collaborative visual affordance for dual-gripper object manipulation. *arXiv preprint arXiv:2207.01971*, 2022. 2