
RH20T: A Comprehensive Robotic Dataset for Learning Diverse Skills in One-Shot

Hao-Shu Fang, Hongjie Fang, Zhenyu Tang, Jirong Liu, Chenxi Wang,
Junbo Wang, Haoyi Zhu and Cewu Lu
Shanghai Jiao Tong University

fhaoshu@gmail.com, {galaxies, tang_zhenyu, jirong, wx1997, sjtuwjb3589635689, zhuhaoyi, lucewu}@sjtu.edu.cn

Abstract

A key challenge for robotic manipulation in open domains is how to acquire diverse and generalizable skills for robots. Recent progress in one-shot imitation learning and robotic foundation models have shown promise in transferring trained policies to new tasks based on demonstrations. This feature is attractive for enabling robots to acquire new skills and improve their manipulative ability. However, due to limitations in the training dataset, the current focus of the community has mainly been on simple cases, such as push or pick-place tasks, relying solely on visual guidance. In reality, there are many complex skills, some of which may even require both visual and tactile perception to solve. This paper aims to unlock the potential for an agent to generalize to hundreds of real-world skills with multi-modal perception. To achieve this, we have collected a dataset comprising over 110,000 *contact-rich* robot manipulation sequences across diverse skills, contexts, robots, and camera viewpoints, all collected *in the real world*. Each sequence in the dataset includes visual, force, audio, and action information. Moreover, we also provide a corresponding human demonstration video and a language description for each robot sequence. We have invested significant efforts in calibrating all the sensors and ensuring a high-quality dataset. The dataset is made publicly available.

1 Introduction

Robotic manipulation requires the robot to control its actuator and change the environment following a task specification. Enabling robots to learn new skills with minimal effort is one of the ultimate goals of the robot learning community. Recent research in one-shot imitation learning [10, 13] and emerging foundation models [3, 5] draw an exciting picture of transferring trained policies to a new task given a demonstration. This paper shares the same aspiration.

While the future is promising, most research in robotics only demonstrates the effectiveness of their algorithms on simple cases, such as pushing, picking, and placing objects in the real world. Two main factors hinder the exploration of more complex tasks in this direction. Firstly, there is a lack of large and diverse robotic manipulation datasets in this field [3], despite the community’s long-standing eagerness for such datasets. The fundamental problem stems from the huge barriers associated with data acquisition. These challenges include the arduous task of configuring diverse robot platforms, creating varied environments, and gathering manipulation trajectories, which require significant effort and resources. Secondly, most methods focus solely on visual guidance control, yet it has been observed in physiology that humans with impaired digital sensibility struggle to accomplish many daily manipulations with visual guidance alone [20]. This indicates that more sensory information should be considered in order to learn various manipulations in open environments.

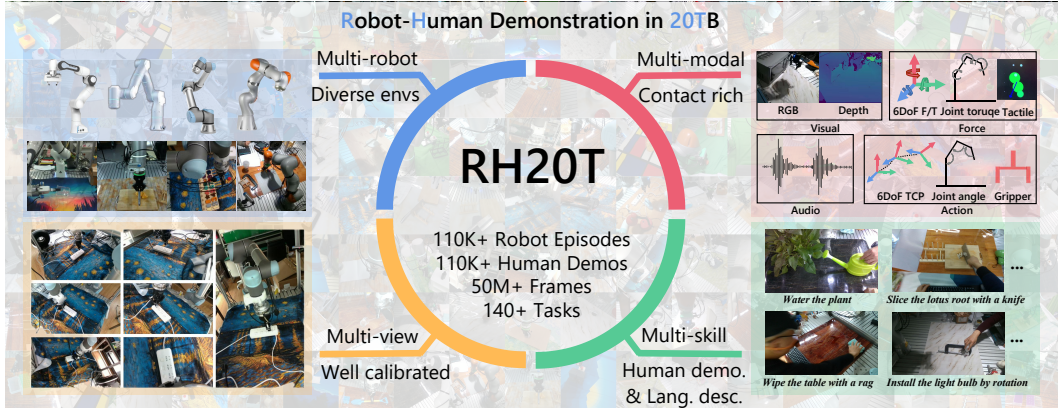


Figure 1: Overview of our RH20T dataset. We adopt multiple robots and setup diverse environments for the data collection. The robot manipulation episodes include multi-modal visual, force, audio and action data. For each episode, we collect the manipulation process with well calibrated multi-view cameras. Our dataset contains diverse robotic manipulation skills and each episode has a corresponding human demonstration and language description. In total, we provide over 110K robot episodes and 110K corresponding human demonstration. The dataset contains over 50 million frames and over 140 tasks.

To address these problems, we revisit the data collection process for robotic manipulation. In most imitation learning literature, expert robot trajectories are manually collected using simplified user interfaces like 3D mice, keyboards, or VR remotes. However, these control methods are inefficient and pose safety risks when the robot engages in rich-contact interactions with the environment. The main reasons are the unintuitive nature of controlling with a 3D mouse or keyboard, and the inaccuracies resulting from motion drifting when using a VR remote. Additionally, tele-operation without force feedback degrades manipulation efficiency for humans. In this paper, we equipped the robot with a force-torque sensor and employed a haptic device with force rendering for precise and efficient data collection. With the goal that the dataset should be representative, generalized, diverse and close to reality, we collect around 150 skills with complicated actions other than simple pick-place. These skills were either selected from RLBench [18] and MetaWorld [40], or proposed by ourselves. Many skills require the robot to engage in contact-rich interactions with the environment, such as cutting, plugging, slicing, pouring, folding, rotating, etc. We have used multiple different robot arms commonly found in labs worldwide to collect our dataset. The diversity in robot configurations can also aid algorithms in generalizing to other robots.

So far, we have collected around 110,000 sequences of robotic manipulation and 110,000 corresponding human demonstration videos for the same skills. This amounts to over 40 million frames of images for the robotic manipulation sequences and over 10 million frames for the human demonstrations. Each robot sequence contains abundant visual, tactile, audio, and proprioception information from multiple sensors. The dataset is carefully organized, and *we believe that a dataset with such diversity and scale is crucial for the future emergence of foundation models in general skill learning*, as promising progress has been witnessed in the NLP and CV communities [6, 31, 22].

2 RH20T Dataset

We introduce our robotic manipulation dataset, Robot-Human demonstration in 20TB (RH20T), to the community. Fig. 1 shows an overview of our dataset.

2.1 Properties of RH20T

RH20T is designed with the objective of enabling general robotic manipulation, which means that the robot can perform various skills based on a task description, typically a human demonstration video, while minimizing the notion of rigid tasks. The following properties are emphasized to fulfill this objective.

Diversity The diversity of RH20T encompasses multiple aspects. To ensure task diversity, we selected 48 tasks from RLBench [18], 29 tasks from MetaWorld [40], and introduced 70 self-proposed tasks that are frequently encountered and achievable by robots. In total, it contains 147

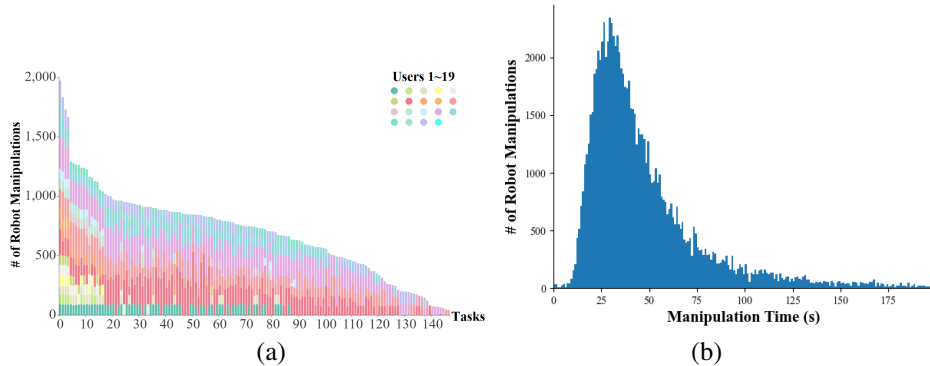


Figure 2: (a) Statistics on the amount of robotic manipulation for different tasks. (b) Statistics on the execution time of different robotic manipulations in our dataset.

tasks, consisting of 42 skills (*i.e.*, verbs). Hundreds of objects were collected to accomplish these tasks. To ensure applicability across different robot configurations, we used 4 popular robot arms, 4 different robotic grippers, and 3 types of force-torque sensors, resulting in 7 robot configurations. Details about the robot configurations are provided in Appendix B.

To enhance environment diversity, we frequently replaced over 50 table covers with different textures and materials, and introduced irrelevant objects to create distractions. Manipulations were performed by tens of volunteers, ensuring diverse trajectories. To increase state diversity, for each skill, volunteers were asked to change the environmental conditions and repeat the manipulation 10 times, including variations in object instances, locations, and more. Additionally, we conducted robotic manipulation experiments involving human interference, both in adversarial and cooperative settings.

Multi-Modal We believe that the future of robotic manipulation lies in multi-modal approaches, particularly in open environments, where data from different sensors will become increasingly accessible with advancements in technology. In the current version of RH20T, we provide visual, tactile, audio, and proprioception information. Visual perception includes RGB, depth, and binocular IR images from three types of cameras. Tactile perception includes 6 DoF force-torque measurements at the robot’s wrist, and some sequences also include fingertip tactile information. Audio data includes recordings from both in-hand and global sources. Proprioception encompasses joint angles/torques, end-effector Cartesian pose and gripper states. All information is collected at the highest frequency supported by our workstation and saved with corresponding timestamps, and the details are given in Appendix B.

Scale Our dataset consists of over 110,000 robot sequences and an equal number of human sequences, with more than 50 million images collected in total. On average, each skill contains approximately 750 robot manipulations. Fig. 2 (a) provides a detailed breakdown of the number of manipulations across different tasks in the dataset, showing a relatively uniform distribution. Fig. 2 (b) presents statistics on the manipulation time for each sequence in our dataset. Most sequences have durations ranging from 10 to 100 seconds. With its substantial volume of data, our dataset stands as the largest in our community at present.

2.2 Data Collection and Processing

Unlike previous methods that simplify the tele-operation interface using 3D mice, VR remotes, or mobile phones, we place emphasis on the importance of intuitive and accurate tele-operation in collecting contact-rich robot manipulation data. Without proper tele-operation, the robot could easily collide with the environment and generate significant forces, triggering emergency stops. Consequently, previous works either avoid contact [19] or operate at reduced speeds to mitigate these risks.

Collection Fig. 3 (a) shows an example of our data collection platform. Each platform contains a robot arm with force-torque sensor, gripper and 1-2 inhand cameras, 8-10 global cameras, 2 microphones, a haptic device, a pedal and a data collection workstation. All the cameras are

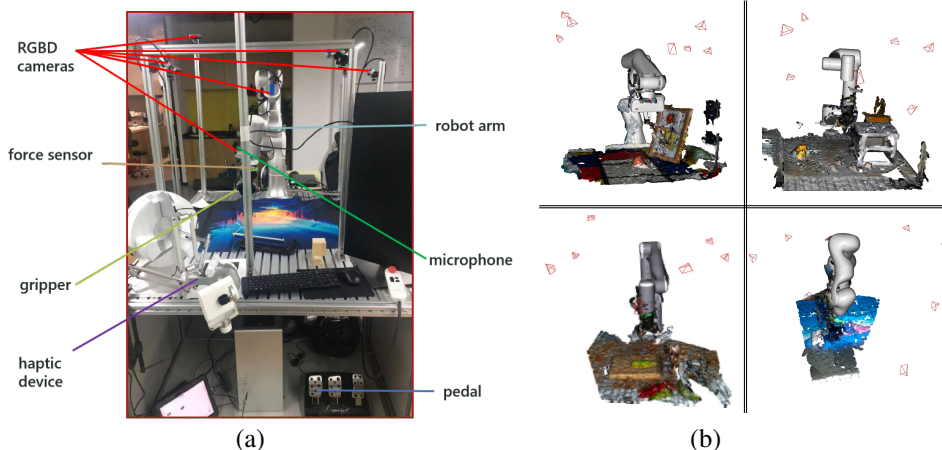


Figure 3: (a) Illustration of our data collection platform. (b) We display the point cloud generated by fusing the RGBD data from the multi-view cameras mounted in our data collection platform. The red pyramids indicate the camera poses.

extrinsically calibrated before conducting the manipulation. The human demonstration video is collected on the same platform by human with an extra ego-centric camera. Tens of volunteers conducted the robotic manipulation according to our task lists and text description. We make our tele-operation pretty intuitive and the average training time is less than 1 hour. The volunteers are also required to specify ending time of the task and give a rating from 0 to 9 after finishing each manipulation. 0 denotes the robot enters the emergency state (e.g., hard collision), 1 denotes the task fails and 2-9 denotes their evaluation of the manipulation quality. The success and failure cases have a ratio of around 10:1 in our dataset.

Processing We preprocess the dataset to provide a coherent data interface. The coordinate frame of all robots and force-torque sensors are aligned. Different force-torque sensors are tared carefully. The end-effector Cartesian pose and the force-torque data are transformed into the coordination system of each camera. Manual validation is performed for each scene to ensure the camera calibration quality. Fig. 3 (b) shows an illustration of rendering different component of the data in a unified coordinate frame and demonstrates the high-quality of our dataset. The detailed data format and data access APIs are provided on our website.

3 Discussion and Conclusion

In this paper we present the RH20T dataset for diverse robotic skill learning. We believe it can facilitate many areas in robotics, especially for robotic manipulation in novel environments. We open source the dataset and hope to promote the development of our community. In the future, we hope to extend our dataset to broader robotic manipulation, including dual-arm and multi-finger dexterous manipulation.

References

- [1] Michal Bednarek, Piotr Kicki, and Krzysztof Walas. On robustness of multi-modal fusion—robotics perspective. *Electronics*, 9(7):1152, 2020.
- [2] Homanga Bharadhwaj, Jay Vakil, Mohit Sharma, Abhinav Gupta, Shubham Tulsiani, and Vikash Kumar. Roboagent: Generalization and efficiency in robot manipulation via semantic augmentations and action chunking. *arXiv preprint arXiv:2309.01918*, 2023.
- [3] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- [4] Alessandro Bonardi, Stephen James, and Andrew J Davison. Learning one-shot imitation from humans without humans. *IEEE Robotics and Automation Letters*, 5(2):3533–3539, 2020.

- [5] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alexander Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil J. Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael S. Ryoo, Grecia Salazar, Pannag R. Sanketi, Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. RT-1: robotics transformer for real-world control at scale. In *Robotics: Science and Systems (RSS)*, 2023.
- [6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems (NeurIPS)*, 33: 1877–1901, 2020.
- [7] Shaowei Cui, Rui Wang, Junhang Wei, Jingyi Hu, and Shuo Wang. Self-attention based visual-tactile fusion learning for predicting grasp outcomes. *IEEE Robotics and Automation Letters*, 5(4):5827–5834, 2020.
- [8] Sudeep Dasari and Abhinav Gupta. Transformers for one-shot imitation learning. In *Conference on Robot Learning (CoRL)*, pages 2071–2084. PMLR, 2020.
- [9] Sudeep Dasari, Frederik Ebert, Stephen Tian, Suraj Nair, Bernadette Bucher, Karl Schmeckpeper, Siddharth Singh, Sergey Levine, and Chelsea Finn. Robonet: Large-scale multi-robot learning. In *Conference on Robot Learning (CoRL)*, volume 100, pages 885–897. PMLR, 2019.
- [10] Yan Duan, Marcin Andrychowicz, Bradley Stadie, OpenAI Jonathan Ho, Jonas Schneider, Ilya Sutskever, Pieter Abbeel, and Wojciech Zaremba. One-shot imitation learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 30, 2017.
- [11] Mark Edmonds, Feng Gao, Xu Xie, Hangxin Liu, Siyuan Qi, Yixin Zhu, Brandon Rothrock, and Song-Chun Zhu. Feeling the force: Integrating force and pose for fluent discovery through imitation learning to open medicine bottles. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3530–3537. IEEE, 2017.
- [12] Nima Fazeli, Miquel Oller, Jiajun Wu, Zheng Wu, Joshua B Tenenbaum, and Alberto Rodriguez. See, feel, act: Hierarchical learning for complex manipulation skills with multisensory fusion. *Science Robotics*, 4(26):eaav3123, 2019.
- [13] Chelsea Finn, Tianhe Yu, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot visual imitation learning via meta-learning. In *Conference on Robot Learning (CoRL)*, pages 357–368. PMLR, 2017.
- [14] Maxwell Forbes, Michael Chung, Maya Cakmak, and Rajesh Rao. Robot programming by demonstration with crowdsourced action fixes. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, volume 2, pages 67–76, 2014.
- [15] De-An Huang, Suraj Nair, Danfei Xu, Yuke Zhu, Animesh Garg, Li Fei-Fei, Silvio Savarese, and Juan Carlos Niebles. Neural task graphs: Generalizing to unseen tasks from a single video demonstration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8565–8574, 2019.
- [16] Tiancheng Huang, Feng Zhao, and Donglin Wang. One-shot imitation learning on heterogeneous associated tasks via conjugate task graph. In *International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2021.
- [17] Stephen James, Michael Bloesch, and Andrew J Davison. Task-embedded control networks for few-shot imitation learning. In *Conference on Robot Learning (CoRL)*, pages 783–795. PMLR, 2018.
- [18] Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J Davison. Rlbench: The robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 5(2): 3019–3026, 2020.
- [19] Eric Jang, Alex Irpan, Mohi Khansari, Daniel Kappler, Frederik Ebert, Corey Lynch, Sergey Levine, and Chelsea Finn. Bc-z: Zero-shot task generalization with robotic imitation learning. In *Conference on Robot Learning (CoRL)*, pages 991–1002. PMLR, 2021.

- [20] Roland S Johansson, J Randall Flanagan, and Roland S Johansson. Sensory control of object manipulation. *Sensorimotor control of grasping: Physiology and pathophysiology*, pages 141–160, 2009.
- [21] Dmitry Kalashnikov, Jacob Varley, Yevgen Chebotar, Benjamin Swanson, Rico Jonschkowski, Chelsea Finn, Sergey Levine, and Karol Hausman. Mt-opt: Continuous multi-task robotic reinforcement learning at scale. *arXiv preprint arXiv:2104.08212*, 2021.
- [22] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv:2304.02643*, 2023.
- [23] Michelle A Lee, Yuke Zhu, Peter Zachares, Matthew Tan, Krishnan Srinivasan, Silvio Savarese, Li Fei-Fei, Animesh Garg, and Jeannette Bohg. Making sense of vision and touch: Learning multimodal representations for contact-rich tasks. *IEEE Transactions on Robotics*, 36(3): 582–596, 2020.
- [24] Fengming Li, Qi Jiang, Wei Quan, Shibo Cai, Rui Song, and Yibin Li. Manipulation skill acquisition for robotic assembly based on multi-modal information description. *IEEE Access*, 8: 6282–6294, 2019.
- [25] Corey Lynch and Pierre Sermanet. Language conditioned imitation learning over unstructured data. In *Robotics: Science and Systems (RSS)*, 2021.
- [26] Zhao Mandi, Fangchen Liu, Kimin Lee, and Pieter Abbeel. Towards more generalizable one-shot visual imitation learning. In *IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2022.
- [27] Ajay Mandlekar, Yuke Zhu, Animesh Garg, Jonathan Booher, Max Spero, Albert Tung, Julian Gao, John Emmons, Anchit Gupta, Emre Orbay, et al. Roboturk: A crowdsourcing platform for robotic skill learning through imitation. In *Conference on Robot Learning (CoRL)*, pages 879–893. PMLR, 2018.
- [28] Peter Pastor, Heiko Hoffmann, Tamim Asfour, and Stefan Schaal. Learning and generalization of motor skills by learning from demonstration. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 763–768. IEEE, 2009.
- [29] Deepak Pathak, Parsa Mahmoudieh, Guanghao Luo, Pulkit Agrawal, Dian Chen, Yide Shentu, Evan Shelhamer, Jitendra Malik, Alexei A Efros, and Trevor Darrell. Zero-shot visual imitation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 2050–2053, 2018.
- [30] Rouhollah Rahmatizadeh, Pooya Abolghasemi, Ladislau Bölöni, and Sergey Levine. Vision-based multi-task manipulation for inexpensive robots using end-to-end learning from demonstration. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 3758–3765. IEEE, 2018.
- [31] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning (ICML)*, pages 8821–8831. PMLR, 2021.
- [32] Nathan Ratliff, J Andrew Bagnell, and Siddhartha S Srinivasa. Imitation learning for locomotion and manipulation. In *2007 7th IEEE-RAS International Conference on Humanoid Robots*, pages 392–397. IEEE, 2007.
- [33] Pratyusha Sharma, Lekha Mohan, Lerrel Pinto, and Abhinav Gupta. Multiple interactions made easy (mime): Large scale demonstrations data for imitation. In *Conference on Robot Learning (CoRL)*, pages 906–915. PMLR, 2018.
- [34] Simon Stepputtis, Joseph Campbell, Mariano Phielipp, Stefan Lee, Chitta Baral, and Heni Ben Amor. Language-conditioned imitation learning for robot manipulation tasks. *Advances in Neural Information Processing Systems (NeurIPS)*, 33:13139–13150, 2020.
- [35] Homer Walke, Kevin Black, Abraham Lee, Moo Jin Kim, Max Du, Chongyi Zheng, Tony Zhao, Philippe Hansen-Estruch, Quan Vuong, Andre He, et al. Bridgedata v2: A dataset for robot learning at scale. *arXiv preprint arXiv:2308.12952*, 2023.
- [36] Zheng Wu, Wenzhao Lian, Vaibhav Unhelkar, Masayoshi Tomizuka, and Stefan Schaal. Learning dense rewards for contact-rich manipulation tasks. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6214–6221. IEEE, 2021.

- [37] Taozheng Yang, Ya Jing, Hongtao Wu, Jiafeng Xu, Kuankuan Sima, Guangzeng Chen, Qie Sima, and Tao Kong. Moma-force: Visual-force imitation for real-world mobile manipulation. *arXiv preprint arXiv:2308.03624*, 2023.
- [38] Sarah Young, Dhiraj Gandhi, Shubham Tulsiani, Abhinav Gupta, Pieter Abbeel, and Lerrel Pinto. Visual imitation made easy. In *Conference on Robot Learning (CoRL)*, volume 155, pages 1992–2005. PMLR, 2020.
- [39] Tianhe Yu, Chelsea Finn, Sudeep Dasari, Annie Xie, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot imitation from observing humans via domain-adaptive meta-learning. In *Robotics: Science and Systems (RSS)*, 2018.
- [40] Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on Robot Learning*, pages 1094–1100. PMLR, 2019.
- [41] Tianhao Zhang, Zoe McCarthy, Owen Jow, Dennis Lee, Xi Chen, Ken Goldberg, and Pieter Abbeel. Deep imitation learning for complex manipulation tasks from virtual reality teleoperation. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 5628–5635. IEEE, 2018.
- [42] Allan Zhou, Eric Jang, Daniel Kappler, Alex Herzog, Mohi Khansari, Paul Wohlhart, Yunfei Bai, Mrinal Kalakrishnan, Sergey Levine, and Chelsea Finn. Watch, try, learn: Meta-learning from demonstrations and rewards. In *International Conference on Learning Representations (ICLR)*, 2019.

Appendices

A Related work

We briefly review related works in robotic manipulation datasets, zero/one-shot imitation learning, and vision-force learning methods.

Dataset Our community has been striving to create a large-scale and representative dataset for a significant period of time. Previous research in one-shot imitation learning has either collected robot manipulation data in the real world [13] or in simulation [26]. However, their datasets are usually small and the tasks are simple. Some attempts have been made to create large-scale real robot manipulation datasets [9, 14, 19, 21, 27, 33]. For example, RoboTurk [27] developed a crowd-sourcing platform and collected data on three tasks using mobile phone-based tele-operation. MIME [33] collected 20 types of manipulations using Baxter with kinesthetic teaching, but they were limited to a single robot and simple environments. RoboNet [9] gathered a significant amount of robot trajectories with various robots, grippers, and environments. However, it mainly consists of random walking episodes due to the challenges of performing meaningful skills. BC-Z [19] presents a manipulation collection of 100 “tasks”, but as pointed out in [26], they are combinations of 9 verbs and 6-15 objects. Similarly, RT-1 [5] and RoboSet [2] also collect large-scale manipulation datasets but focus on a limited set of skills. Concurrently to our work, BridgeData V2 [35] collects a dataset with 13 skills across 24 environments. In this paper, we present a larger dataset with a wider range of skills and environments, with more comprehensive information. More importantly, all previous datasets put less emphasize on contact-rich manipulation. Our dataset focus more in this case and include the crucial force modality during manipulation.

Zero/One-shot Imitation Learning The objective of training policies that can transfer to new tasks based on robot/human demonstrations is not new. Early works [32, 28, 14] focused on imitation learning using high-level states such as trajectories. Recently, researchers [13, 10, 41, 17, 39, 30, 29, 42, 15, 34, 4, 38, 8, 25, 19, 26] have started exploring raw-pixel inputs with the advancement of deep neural networks. Additionally, the requirement of demonstrations has been reduced by eliminating the need for actions. Recent approaches have explored various one-shot task descriptors, including images [17, 4], language [34, 25, 5, 2], robot video [13, 8, 26], or human video [39, 19]. These methods can be broadly classified into three categories: model-agnostic meta-learning [13, 39, 17, 4, 42], conditional behavior cloning [10, 8, 19, 5, 26], and task graph construction [15, 16]. While significant progress has been made in this direction, these approaches only consider visual observations and primarily focus on simple robotic manipulations such as reach, pick, push, or place. Our dataset offers the opportunity to take a step further by enabling the learning of *hundreds* of skills that require *multi-modal perception* within a single imitation learning model.

Multi-Modal Learning of Vision and Force Force perception plays a crucial role in manipulation tasks, providing valuable and complementary information when visual perception is occluded. The joint modeling of vision and force in robotic manipulation has recently garnered interest within the research community [11, 24, 12, 23, 1, 7, 36]. However, most of these studies overlook the asynchronous nature of different modalities and simply concatenate the signals before or after the neural network. Moreover, the existing research primarily focuses on designing multi-modal learning algorithms for specific tasks, such as grasping [7], insertion [23], twisting [11], or playing Jenga [12]. A recent attempt [37] explores jointly imitating the action and wrench on 6 tasks respectively. Overall, the question of how to effectively handle multi-modal perception at different frequencies for various skills in a coherent manner remains open in robotics. Our dataset presents an opportunity for exploring multi-sensory learning across diverse real-world skills.

B Data details

Conf.	Robot	Gripper	6DoF F/T Sensor	Tactile
Cfg 1	Flexiv	Dahuan AG95	OptoForce	N/A
Cfg 2	Flexiv	Dahuan AG95	ATI Axia80-M20	N/A
Cfg 3	UR5	WSG50	ATI Axia80-M20	N/A
Cfg 4	UR5	Robotiq-85	ATI Axia80-M20	N/A
Cfg 5	Franka	Franka	Franka	N/A
Cfg 6	Kuka	Robotiq-85	ATI Axia80-M20	N/A
Cfg 7	Kuka	Robotiq-85	ATI Axia80-M20	uSkin

Table 1: Hardware specification of different configurations.

Conf.	Modal	Size	Frequency
Cfg 1-7	RGB image	1280×720×3	10 Hz
	Depth image	1280×720	10 Hz
	Binocular IR image	1280×720	10 Hz
	Robot joint angle	6 / 7	10 Hz
	Robot joint torque	6 / 7	10 Hz
	Gripper Cartesian pose	6 / 7	100 Hz
	Gripper width	1	10 Hz
	6DoF F/T	6	100 Hz
	Audio	N/A	30 Hz
Cfg 7	Tactile	2×16×3	200 Hz

Table 2: Data information of different configurations. The first 9 data modality are the same for all robot configurations. The last data modality of fingertip tactile sensing is only available in Cfg 7.