

RoboCrowd: Scaling Robot Data Collection through Crowdsourcing

Suvir Mirchandani, David D. Yuan, Kaylee Burns, Md Sazzad Islam,
Tony Z. Zhao, Chelsea Finn, Dorsa Sadigh
Stanford University
<https://robocrowd.github.io>

Abstract: Imitation learning from large-scale human demonstrations has emerged as a promising paradigm for training robot policies. However, collecting demonstrations is burdensome for expert operators. We introduce a new data collection paradigm, RoboCrowd, which distributes the workload by utilizing crowdsourcing principles and incentive design. We build RoboCrowd on top of ALOHA [1]—a bimanual platform that supports data collection via puppeteering—to explore the design space for crowdsourcing in-person demonstrations in a public environment. We propose three classes of incentive mechanisms to appeal to users’ varying sources of motivation for interacting with the system: material rewards, intrinsic interest, and social comparison. We instantiate these incentives through tasks that include physical rewards, engaging or challenging manipulations, as well as gamification elements such as a leaderboard. We conduct a large-scale, two-week field experiment in which the platform is situated in a university café. Over 200 individuals independently volunteered to provide a total of over 800 interaction episodes. Our findings validate the proposed incentives as mechanisms for shaping users’ data quantity and quality. Further, we demonstrate that the crowdsourced data can serve as useful pre-training data for policies fine-tuned on expert demonstrations—boosting performance up to 20% compared to when this data is not available. These results suggest the potential for RoboCrowd to reduce the burden of robot data collection by carefully implementing crowdsourcing and incentive design principles.

1 Introduction

Imitation learning (IL) has become a popular paradigm for training robot policies [1–5]. However, modern IL algorithms continue to have significant data requirements especially as tasks increase in number and variety—on the order of hundreds to thousands of demonstrations [5, 6]. Prior efforts to scale up real-world data collection include pooling demonstration data across different institutions [6–8], which has amortized the cost of real-robot data collection to a degree. A fundamental limitation to scaling up is that the source of demonstrations is primarily researchers or designated operators. To explore ways to scale up robot data collection, we ask: *Who* can effectively collect robot data, and *how* might they be incentivized to do so?

To tackle this problem, we look to a large body of work outside of robotics which studies strategies for incentivizing people in crowdsourced data labeling tasks [9–14]. The goal of these works is to align the incentives of crowdworkers with researchers’ goals of labeling a given dataset—for example, *gamifying* the data labeling process [11]. Our key idea is to build a system that leverages similar ideas for robot data collection—i.e., *aligning human incentives to provide robot demonstration data*. We propose **RoboCrowd**, a framework for incentive design in the context of crowdsourced robot data collection. Our framework centers five key properties: public accessibility, capability, intuitiveness, safety, and gamification. We incorporate three classes of incentives to appeal to users’ varying sources of motivation for interacting with the system. To instantiate the framework, we build upon ALOHA [1]—a bimanual platform for robot teleoperation. We deploy the system in a field experiment in which the robot is situated near a university café, where users participate in a self-guided, gamified data collection experience. Over 200 individuals independently volunteered to provide a total of over 800 interaction episodes. We compile the crowdsourced interactions into a dataset and annotate each trajectory with quality scores and task labels. We additionally validate the incentive mechanisms for shaping user interactions with the robot. Finally, we analyze the usefulness of the crowdsourced data for

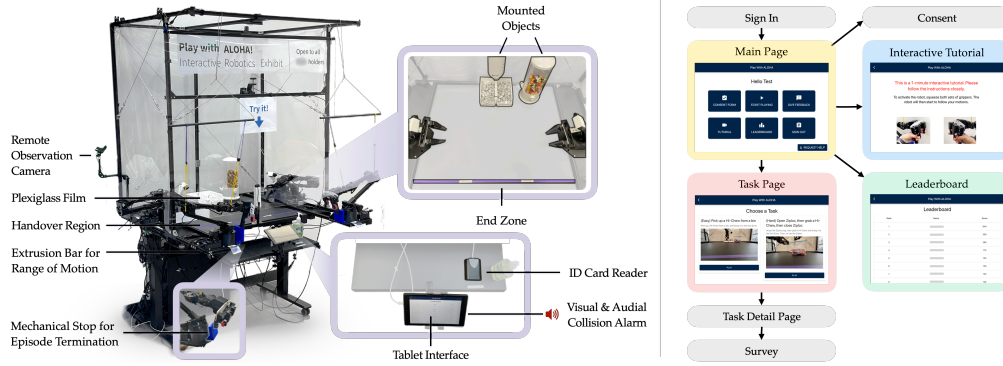


Figure 1: **System Overview.** (Left) RoboCrowd uses the ALOHA robot [1]. Users can perform tasks in scenes put in place by the scene designer; tasks may include physical rewards that the user can bring to the End Zone and access via the Handover Region. (Right) Users are guided by a GUI on a tablet. Functionalities include an Interactive Tutorial, a Task Page, and a Leaderboard. For additional details, please see Appendix E.

training policies. We demonstrate that the crowdsourced data can serve as useful pre-training data when fine-tuning on expert demonstrations, boosting policy performance up to 20% compared to expert-only policies.

2 Related Work

Crowdsourcing is a well-studied technique in human-computer interaction, often used for collecting data labels from a large set of users, with a variety of applications from computer vision to natural language processing [14–23]. While many works utilize platforms such as Amazon Mechanical Turk [24] and Prolific [25] to pay crowdworkers for data labels, other works consider how to *incentivize* crowdworkers via other incentives beyond direct payment to gather data [12, 26–29]. Crowdsourcing has also been an attractive approach for collecting data in robotics in recent years. Prior works have attempted to crowdsource robot data via remote teleoperation in simulation or via web interfaces. RoboTurk [30, 31] develops a smartphone interface to allow crowdworkers on Mechanical Turk to collect demonstrations remotely, and shows the potential of using crowdsourced data to aid policy learning. Several works have developed new interfaces to make robot demonstration collection more distributed. Recent works [32–34] design new hardware interfaces—e.g., sensorized hand-held grippers or portable motion capture systems—to allow for demonstration collection in the real-world without needing access to a physical robot. However, crowdsourcing data with these interfaces is not immediately possible since it still requires data collectors to have access to this custom hardware. In this work, we leverage an existing interface (puppeteering via ALOHA [1], which enables precise bimanual manipulation at a low-cost) and choose to situate it directly in a public space to make it accessible to data collectors. To make scaling up data collection possible, we design the system so it can be used by non-experts.

3 RoboCrowd

We apply incentive design to the collection of robot demonstrations for imitation learning, and develop a system to collect demonstrations directly from the public. We establish a set of desired properties for our system to enable crowdsourcing robot data: [P1] *publicly accessible*, [P2] *capable hardware*, and [P3] *intuitive* and [P4] *safe* for novices, and [P5] *gamified*. Additionally, we design incentive mechanisms to shape these interactions into useful data. We expect that crowdworkers vary in their motivations; we therefore design for three incentive mechanisms: [M1] *material rewards* (e.g., physical rewards for completing a task), [M2] *intrinsic interest* (e.g., challenging or engaging tasks), and [M3] *social comparison* (e.g., a leaderboard). This section explains how we meet these desiderata through our hardware and software design.

Hardware Design. We select ALOHA [1], a system for bimanual teleoperation, as the base platform for our system. ALOHA consists of two “follower” arms (ViperX) that are controlled via puppeteering with two “leader” arms (WidowX). We choose to use the ALOHA platform due to its low-cost, reparability, as well as its ability for collecting data for a wide task range. Fig. 1 illustrates a set of enhancements to outfit ALOHA for public use to achieve our desired properties and enable crowdsourcing. First, we implement mechanisms for user and robot safety (P4): (a) collision avoidance to prevent self-collisions, achieved via a parallel MuJoCo [35] simulator, as well as a visual-audial alarm when the robot is near collision; (b) plexiglass and



Figure 2: **Scenes.** BinScene, Bin+DispenserScene, and Bin+ZiplocScene, and the objects relevant to the tasks (hi-chew, tootsie-roll, hershey-kiss, jelly-bean, hi-chew-bin, hi-chew-ziploc).

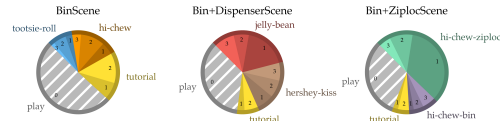


Figure 3: **Dataset composition by number of time steps for each of our three scenes.** Tasks receive quality scores from 1 to 3 (higher is better) which are also indicated by brighter shades. Tutorial data receives a score of 1 or 2. Play data always receives a score of 0.

vinyl film to cover all sides of the ALOHA workcell to enclose the puppet arms; (c) extended extrusion bars on the leader arms to increase the distance between users and leader arms; (d) mounting of scene props (such as bins and dispensers) to mitigate scene damage; and (e) a remote observation camera for the scene designer to periodically monitor the scene. We also include enhancements to increase the intuitiveness of the platform for members of the public (P3): (a) a tablet interface, described in the next section; (b) a mechanical stop for users to automatically terminate episodes by resting the puppet arms. To enable a gamified setup (P5), we utilize (a) an ID card reader to authenticate and track users and (b) demarcate an “End Zone” within scenes, where a user can place physical rewards and access them via a handover region at the bottom of the plexiglass casing. Given its ability to perform versatile tasks, ALOHA satisfies our capability goal (P2). We physically situate it in a public environment (Section 4) to make it accessible to crowd users (P1).

Software Design. To make operating the robot intuitive (P2) for members of the public, we implement a tablet application to complement the hardware platform and guide users through the operation process (Fig. 1; right). The interface additionally features a variety of elements of gamification (P5). We develop an onboarding process for new users to sign-in and receive a tutorial to familiarize themselves with the platform. We design our onboarding process to be efficient and interactive: users begin by tapping their university ID card on a card reader, which directs them to a Sign In page to create a *user profile*. Users are then directed to complete a consent form and an interactive tutorial to learn how to puppeteer the robot (Fig. 1; right). The tutorial contains four steps and takes less than one minute to complete. We detail the stages of the interactive tutorial in Appendix E. After completing the tutorial, users can choose to enter a *Task Page* where they see videos of different tasks they can complete in the scene (Fig. 1; right). In service of P5, we use gamified verbiage and elements throughout the interface (e.g. a *Start Playing* button, and a *countdown timer* on performing tasks). Specifically for M3, we implement a point system where users receive points for completing tasks, which are tallied and visible on a *Leaderboard Page*, where users can see how their scores rank compared to other users (Fig. 1; right). We describe implementation details of the software architecture in Appendix E.

4 Experiments

We utilize RoboCrowd to collect a crowdsourced dataset over a two-week period in a public university café. We instantiate three types of incentive mechanisms (M1-M3) to appeal to users’ varying motivations, and design scenes in order to verify if these mechanisms can shape demonstration quantity and quality.

Scene Design. On each day of crowdsourcing, two of six tasks are made available to users, with different pairs corresponding to different scenes (Fig. 2). BinScene contains bins with two types of candies for single arm bin-picking tasks (hi-chew and tootsie-roll). Bin+DispenserScene contains the same bins with a single type of candy (hershey-kiss), as well as a cup dispenser and a jelly bean dispenser (jelly-bean). Bin+ZiplocScene contains the same bins with a single candy type (hi-chew-bin) as well as a closed Ziploc bag full of candies (hi-chew-ziploc). Please see Appendix A for task details.

We observe significant engagement with RoboCrowd over the two-week collection period: there were $N = 231$ unique users in total. We collect 129 interaction episodes in BinScene (Day 1), 381 in Bin+DispenserScene (Days 2–5), and 307 in Bin+ZiplocScene (Days 6–11). In aggregate, users spent 54.2% of interaction time performing the preset tasks in the scene, 9.6% on the interactive tutorial, and 36.1% on free-play. In Fig. 3, we show the distribution of tasks and qualities over timesteps for each scene. Qualities are determined on a scale from 1–3 for task-relevant data and a scale of 1–2 for tutorial data based on the smoothness of the user’s motion and whether there is retrying behavior or extraneous movements.

We detail the quality annotation rules in Appendix E, and illustrate sample trajectories in Appendix A.

Effects of Incentives on Data Quantity and Quality. *Material Rewards.* While BinScene contains two bin-picking tasks with nearly identical difficulty, users in aggregate spend $2\times$ as many timesteps performing hi-chew compared to tootsie-roll. This suggests that users devote more interaction time to tasks where the direct material incentive is more preferred (users generally express preferences for Hi-Chews per our offline study). Users also spend a significant amount of time (50.7%) on free-play with the system in BinScene, engaging in behaviors such as trying out more challenging tasks (e.g., attempting to unwrap the candies; see Appendix B). Thus, while material incentives can influence user demonstrations, drivers of intrinsic motivation such as task difficulty also play a role.

Intrinsic Motivation. In Bin+DispenserScene, which contains a harder bin-picking task than in Scene A (hershey-kiss) and a challenging candy dispensing task (jelly-bean), users spend only 35.3% of the time in free-play. Despite the fact that users do not generally prefer Jelly Beans over Hershey Kisses as a material reward, they still spend more ($1.5\times$) time performing the jelly-bean task. This suggests that intrinsic interest can influence users to allocate more time doing harder task compared to easier ones, or engaging in free-play. To probe this effect even when controlling for material reward, we consider Bin+ZiplocScene. Here, the incentive is contained within a closed Ziploc bag which must be opened. The same incentive is available in the bin to be picked. Users spend $4.18\times$ as many timesteps on hi-chew-ziploc compared to hi-chew-bin, again suggesting that intrinsic motivation influences which tasks users perform in the scene.

Social Comparison. To examine how different people respond differently to explicit comparison mechanisms in the system, we record which users visit the Leaderboard Page, and conduct a Mann-Whitney U-test to compare the quantity and quality of demonstrations provided by Leaderboard visitors compared to other users. Fig. 4 illustrates the distribution of quantity (number of interactions) and quality (mean quality score) conditioned on Leaderboard visitation. We find that that visitors of the Leaderboard provide significantly more demonstrations ($p < 0.001$) that are higher quality on average ($p < 0.05$).

Policy Learning. Finally, we study how useful the crowdsourced data is for downstream policy learning. To complement the crowdsourced data, we collect a set of high-quality expert demonstrations for each task: 30 demonstrations for each of hi-chew and tootsie-roll, 60 for hershey-kiss, 80 for hi-chew-bin, and 100 for each of jelly-bean and hi-chew-ziploc. In Table 1, we compare different methods of mixing crowdsourced data and expert data on our six tasks. All policies use ACT [1] with default hyperparameters. Training exclusively with the expert data on each task constitutes the *Expert* setting. *Co-train* refers to naively mixing data from a crowdsourced task (i.e., task-relevant data of any quality) with the expert data. We also compare to *Fine-tune*, which trains in two stages: first co-training on the crowd data and expert data and then fine-tuning on expert data only; for fair comparison, note that *Fine-tune* is trained for fewer total steps (150K) than both *Expert* and *Co-train* (200K). Crowdsourced data provides performance improvements in multiple cases, but the specific effects vary by task. For example, crowdsourced data for the bin-picking tasks can involve low-quality behaviors (i.e., regrasping behavior or grasping multiple items at a time), which may cause the *Co-train* to perform worse than *Expert*, but still provide a useful initialization for *Fine-tune*. We provide additional qualitative analysis of the trained policies in Appendix D.2.

Discussion. We propose and validate a new paradigm for robot data collection via crowdsourcing and incentive design. Crowdsourcing can reduce data collection effort of individual researchers, but also presents challenges of data quality and heterogeneity. Future work can seek to understand the style of different operators and the most effective ways to leverage crowdsourced data during downstream policy learning.

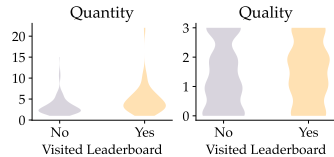


Figure 4: **Quantity and quality by leaderboard use.** Violin plot showing the distribution of quantity and quality of demonstrations for users who did and did not visit the leaderboard.

Task	Expert	Co-train	Fine-tune
hi-chew	37.5%	27.5%	42.5%
tootsie-roll	42.5%	25%	40%
hershey-kiss	20%	32.5%	35%
hi-chew-bin	20%	12.5%	40%
jelly-bean	48.9 ± 18.6	8.9 ± 10.1	19.7 ± 29.7
hi-chew-ziploc	5.4 ± 12.2	17.1 ± 15.8	22.1 ± 14.3

Table 1: **Policy Performance.** Performance of policies trained on expert demonstrations (# Exp.), co-trained on crowd data, and pre-trained on expert+crowd data then fine-tuned on expert data. We conduct 40 trials for each cell. For the long-horizon tasks (jelly-bean, hi-chew-ziploc), we provide a normalized return (out of 100) rather than success rate (see Appendix D for details).

References

- [1] T. Zhao, V. Kumar, S. Levine, and C. Finn. Learning Fine-Grained Bimanual Manipulation with Low-Cost Hardware. In *Proceedings of Robotics: Science and Systems (RSS)*, 2023.
- [2] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, T. Jackson, S. Jesmonth, N. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, I. Leal, K.-H. Lee, S. Levine, Y. Lu, U. Malla, D. Manjunath, I. Mordatch, O. Nachum, C. Parada, J. Peralta, E. Perez, K. Pertsch, J. Quiambao, K. Rao, M. Ryoo, G. Salazar, P. Sanketi, K. Sayed, J. Singh, S. Sontakke, A. Stone, C. Tan, H. Tran, V. Vanhoucke, S. Vega, Q. Vuong, F. Xia, T. Xiao, P. Xu, S. Xu, T. Yu, and B. Zitkovich. RT-1: Robotics Transformer for Real-World Control at Scale. In *arXiv*, 2022.
- [3] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess, A. Dubey, C. Finn, P. Florence, C. Fu, M. G. Arenas, K. Gopalakrishnan, K. Han, K. Hausman, A. Herzog, J. Hsu, B. Ichter, A. Irpan, N. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, I. Leal, L. Lee, T.-W. E. Lee, S. Levine, Y. Lu, H. Michalewski, I. Mordatch, K. Pertsch, K. Rao, K. Reymann, M. Ryoo, G. Salazar, P. Sanketi, P. Sermanet, J. Singh, A. Singh, R. Soricut, H. Tran, V. Vanhoucke, Q. Vuong, A. Wahid, S. Welker, P. Wohlhart, J. Wu, F. Xia, T. Xiao, P. Xu, S. Xu, T. Yu, and B. Zitkovich. RT-2: Vision-Language-Action Models Transfer Web Knowledge to Robotic Control. In *arXiv*, 2023.
- [4] Octo Model Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna, C. Xu, J. Luo, T. Kreiman, Y. Tan, L. Y. Chen, P. Sanketi, Q. Vuong, T. Xiao, D. Sadigh, C. Finn, and S. Levine. Octo: An Open-Source Generalist Robot Policy. In *Proceedings of Robotics: Science and Systems (RSS)*, 2024.
- [5] M. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, R. Rafailov, E. Foster, G. Lam, P. Sanketi, Q. Vuong, T. Kollar, B. Burchfiel, R. Tedrake, D. Sadigh, S. Levine, P. Liang, and C. Finn. OpenVLA: An Open-Source Vision-Language-Action Model. In *Conference on Robot Learning (CoRL)*, 2024.
- [6] O. X.-E. Collaboration, A. O’Neill, A. Rehman, A. Maddukuri, A. Gupta, A. Padalkar, A. Lee, A. Pooley, A. Gupta, A. Mandlekar, A. Jain, A. Tung, A. Bewley, A. Herzog, A. Irpan, A. Khazatsky, A. Rai, A. Gupta, A. Wang, A. Kolobov, A. Singh, A. Garg, A. Kembhavi, A. Xie, A. Brohan, A. Raffin, A. Sharma, A. Yavary, A. Jain, A. Balakrishna, A. Wahid, B. Burgess-Limerick, B. Kim, B. Schölkopf, B. Wulfe, B. Ichter, C. Lu, C. Xu, C. Le, C. Finn, C. Wang, C. Xu, C. Chi, C. Huang, C. Chan, C. Agia, C. Pan, C. Fu, C. Devin, D. Xu, D. Morton, D. Driess, D. Chen, D. Pathak, D. Shah, D. Büchler, D. Jayaraman, D. Kalashnikov, D. Sadigh, E. Johns, E. Foster, F. Liu, F. Ceola, F. Xia, F. Zhao, F. V. Frujeri, F. Stulp, G. Zhou, G. S. Sukhatme, G. Salhotra, G. Yan, G. Feng, G. Schiavi, G. Berseth, G. Kahn, G. Yang, G. Wang, H. Su, H.-S. Fang, H. Shi, H. Bao, H. B. Amor, H. I. Christensen, H. Furuta, H. Walke, H. Fang, H. Ha, I. Mordatch, I. Radosavovic, I. Leal, J. Liang, J. Abou-Chakra, J. Kim, J. Drake, J. Peters, J. Schneider, J. Hsu, J. Bohg, J. Bingham, J. Wu, J. Gao, J. Hu, J. Wu, J. Wu, J. Sun, J. Luo, J. Gu, J. Tan, J. Oh, J. Wu, J. Lu, J. Yang, J. Malik, J. Silvério, J. Hejna, J. Booher, J. Tompson, J. Yang, J. Salvador, J. J. Lim, J. Han, K. Wang, K. Rao, K. Pertsch, K. Hausman, K. Go, K. Gopalakrishnan, K. Goldberg, K. Byrne, K. Oslund, K. Kawaharazuka, K. Black, K. Lin, K. Zhang, K. Ehsani, K. Lekkala, K. Ellis, K. Rana, K. Srinivasan, K. Fang, K. P. Singh, K.-H. Zeng, K. Hatch, K. Hsu, L. Itti, L. Y. Chen, L. Pinto, L. Fei-Fei, L. Tan, L. J. Fan, L. Ott, L. Lee, L. Weihs, M. Chen, M. Lepert, M. Memmel, M. Tomizuka, M. Itkina, M. G. Castro, M. Spero, M. Du, M. Ahn, M. C. Yip, M. Zhang, M. Ding, M. Heo, M. K. Srirama, M. Sharma, M. J. Kim, N. Kanazawa, N. Hansen, N. Heess, N. J. Joshi, N. Suenderhauf, N. Liu, N. D. Palo, N. M. M. Shafiullah, O. Mees, O. Kroemer, O. Bastani, P. R. Sanketi, P. T. Miller, P. Yin, P. Wohlhart, P. Xu, P. D. Fagan, P. Mitrano, P. Sermanet, P. Abbeel, P. Sundaresan, Q. Chen, Q. Vuong, R. Rafailov, R. Tian, R. Doshi, R. Mart’ in-Mart’ in, R. Bajjal, R. Scalise, R. Hendrix, R. Lin, R. Qian, R. Zhang, R. Mendonca, R. Shah, R. Hoque, R. Julian, S. Bustamante, S. Kirmani, S. Levine, S. Lin, S. Moore, S. Bahl, S. Dass, S. Sonawani, S. Song, S. Xu, S. Haldar, S. Karamcheti, S. Adebola, S. Guist,

- S. Nasiriany, S. Schaal, S. Welker, S. Tian, S. Ramamoorthy, S. Dasari, S. Belkhale, S. Park, S. Nair, S. Mirchandani, T. Osa, T. Gupta, T. Harada, T. Matsushima, T. Xiao, T. Kollar, T. Yu, T. Ding, T. Davchev, T. Z. Zhao, T. Armstrong, T. Darrell, T. Chung, V. Jain, V. Vanhoucke, W. Zhan, W. Zhou, W. Burgard, X. Chen, X. Chen, X. Wang, X. Zhu, X. Geng, X. Liu, X. Liangwei, X. Li, Y. Pang, Y. Lu, Y. J. Ma, Y. Kim, Y. Chebotar, Y. Zhou, Y. Zhu, Y. Wu, Y. Xu, Y. Wang, Y. Bisk, Y. Dou, Y. Cho, Y. Lee, Y. Cui, Y. Cao, Y.-H. Wu, Y. Tang, Y. Zhu, Y. Zhang, Y. Jiang, Y. Li, Y. Li, Y. Iwasawa, Y. Matsuo, Z. Ma, Z. Xu, Z. J. Cui, Z. Zhang, Z. Fu, and Z. Lin. Open X-Embodiment: Robotic Learning Datasets and RT-X Models. In *International Conference on Robotics and Automation (ICRA)*, 2024.
- [7] A. Khazatsky, K. Pertsch, S. Nair, A. Balakrishna, S. Dasari, S. Karamcheti, S. Nasiriany, M. K. Srirama, L. Y. Chen, K. Ellis, P. D. Fagan, J. Hejna, M. Itkina, M. Lepert, Y. J. Ma, P. T. Miller, J. Wu, S. Belkhale, S. Dass, H. Ha, A. Jain, A. Lee, Y. Lee, M. Memmel, S. Park, I. Radosavovic, K. Wang, A. Zhan, K. Black, C. Chi, K. B. Hatch, S. Lin, J. Lu, J. Mercat, A. Rehman, P. R. Sanketi, A. Sharma, C. Simpson, Q. Vuong, H. R. Walke, B. Wulfe, T. Xiao, J. H. Yang, A. Yavary, T. Z. Zhao, C. Agia, R. Bajjal, M. G. Castro, D. Chen, Q. Chen, T. Chung, J. Drake, E. P. Foster, J. Gao, D. A. Herrera, M. Heo, K. Hsu, J. Hu, D. Jackson, C. Le, Y. Li, K. Lin, R. Lin, Z. Ma, A. Maddukuri, S. Mirchandani, D. Morton, T. Nguyen, A. O’Neill, R. Scalise, D. Seale, V. Son, S. Tian, E. Tran, A. E. Wang, Y. Wu, A. Xie, J. Yang, P. Yin, Y. Zhang, O. Bastani, G. Berseth, J. Bohg, K. Goldberg, A. Gupta, A. Gupta, D. Jayaraman, J. J. Lim, J. Malik, R. Martín-Martín, S. Ramamoorthy, D. Sadigh, S. Song, J. Wu, M. C. Yip, Y. Zhu, T. Kollar, S. Levine, and C. Finn. DROID: A Large-Scale In-The-Wild Robot Manipulation Dataset. In *Proceedings of Robotics: Science and Systems (RSS)*, 2024.
- [8] S. Dasari, F. Ebert, S. Tian, S. Nair, B. Bucher, K. Schmeckpeper, S. Singh, S. Levine, and C. Finn. RoboNet: Large-Scale Multi-Robot Learning. In *Conference on Robot Learning (CoRL)*, 2019.
- [9] R. Snow, B. O’Connor, D. Jurafsky, and A. Y. Ng. Cheap and Fast - But is it Good? Evaluating Non-Expert Annotations for Natural Language Tasks. In *Conference on Empirical Methods in Natural Language Processing*, 2008.
- [10] A. Sorokin and D. A. Forsyth. Utility data annotation with Amazon Mechanical Turk. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2008.
- [11] L. von Ahn and L. Dabbish. Labeling images with a computer game. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, 2004.
- [12] M. S. Bernstein, D. S. Tan, G. Smith, M. Czerwinski, and E. Horvitz. Collabio: a game for annotating people within social networks. In *ACM Symposium on User Interface Software and Technology (UIST)*, 2009.
- [13] J. Park, R. Krishna, P. Khadpe, L. Fei-Fei, and M. S. Bernstein. AI-Based Request Augmentation to Increase Crowdsourcing Participation. In *Proceedings of the Seventh AAAI Conference on Human Computation and Crowdsourcing, HCOMP*, 2019.
- [14] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, M. S. Bernstein, and L. Fei-Fei. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations. *International Journal of Computer Vision*, 2016.
- [15] N. Zhou, Z. D. Siegel, S. Zarecor, N. Lee, D. A. Campbell, C. M. Andorf, D. Nettleton, C. J. Lawrence-Dill, B. Ganapathysubramanian, J. W. Kelly, and I. Friedberg. Crowdsourcing image analysis for plant phenomics to generate ground truth data for machine learning. *PLOS Computational Biology*, 07 2018.
- [16] S. Ørting, A. Doyle, A. van Hilten, M. Hirth, O. Inel, C. R. Madan, P. Mavridis, H. Spiers, and V. Cheplygina. A Survey of Crowdsourcing in Medical Image Analysis. *Human Computation*, 2019.

- [17] T. W. Cenggoro, F. Tanzil, A. H. Aslamiah, E. K. Karuppiah, and B. Pardamean. Crowdsourcing annotation system of object counting dataset for deep learning algorithm. *IOP Conference Series: Earth and Environmental Science*, 2018.
- [18] M. van Vliet, E. C. Groen, F. Dalpiaz, and S. Brinkkemper. Identifying and Classifying User Requirements in Online Feedback via Crowdsourcing. In *Requirements Engineering: Foundation for Software Quality*, 2020.
- [19] B. Shmueli, J. Fell, S. Ray, and L.-W. Ku. Beyond Fair Pay: Ethical Implications of NLP Crowdsourcing. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, 2021.
- [20] N. Nangia, S. Sugawara, H. Trivedi, A. Warstadt, C. Vania, and S. R. Bowman. What Ingredients Make for an Effective Crowdsourcing Protocol for Difficult NLU Data Collection Tasks? *arXiv*, 2021.
- [21] S. Mishra, D. Khashabi, C. Baral, and H. Hajishirzi. Cross-Task Generalization via Natural Language Crowdsourcing Instructions. *arXiv*, 2022.
- [22] S. Lim, A. Jatowt, M. Färber, and M. Yoshikawa. Annotating and Analyzing Biased Sentences in News Articles using Crowdsourcing. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, May 2020.
- [23] L. B. Chilton, G. Little, D. Edge, D. S. Weld, and J. A. Landay. Cascade: crowdsourcing taxonomy creation. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, 2013.
- [24] K. Crowston. Amazon Mechanical Turk: A Research Tool for Organizations and Information Systems Scholars. In *Shaping the Future of ICT Research. Methods and Approaches*, 2012.
- [25] S. Palan and C. Schitter. Prolific. ac—A subject pool for online experiments. *Journal of Behavioral and Experimental Finance*, 2018.
- [26] L. Von Ahn, B. Maurer, C. McMillen, D. Abraham, and M. Blum. recaptcha: Human-based character recognition via web security measures. *Science*, 2008.
- [27] E. Law and L. Von Ahn. *Human computation*. Morgan & Claypool Publishers, 2011.
- [28] H. Garcia-Molina, M. Joglekar, A. Marcus, A. Parameswaran, and V. Verroios. Challenges in Data Crowdsourcing. *IEEE Transactions on Knowledge and Data Engineering*, 2016.
- [29] J. Huynh, J. Bigham, and M. Eskenazi. A Survey of NLP-Related Crowdsourcing HITs: what works and what does not. *arXiv*, 2021.
- [30] A. Mandlekar, Y. Zhu, A. Garg, J. Booher, M. Spero, A. Tung, J. Gao, J. Emmons, A. Gupta, E. Orbay, S. Savarese, and L. Fei-Fei. RoboTurk: A Crowdsourcing Platform for Robotic Skill Learning through Imitation. In *Conference on Robot Learning*, 2018.
- [31] A. Mandlekar, J. Booher, M. Spero, A. Tung, A. Gupta, Y. Zhu, A. Garg, S. Savarese, and L. Fei-Fei. Scaling robot supervision to hundreds of hours with roboturk: Robotic manipulation dataset through human reasoning and dexterity. In *International Conference on Intelligent Robots and Systems (IROS)*, 2019.
- [32] S. Song, A. Zeng, J. Lee, and T. Funkhouser. Grasping in the wild: Learning 6dof closed-loop grasping from low-cost demonstrations. *IEEE Robotics and Automation Letters (RA-L)*, 2020.
- [33] C. Wang, H. Shi, W. Wang, R. Zhang, L. Fei-Fei, and C. K. Liu. DexCap: Scalable and Portable Mocap Data Collection System for Dexterous Manipulation. In *Proceedings of Robotics: Science and Systems (RSS)*, 2024.

- [34] C. Chi, Z. Xu, C. Pan, E. A. Cousineau, B. Burchfiel, S. Feng, R. Tedrake, and S. Song. Universal Manipulation Interface: In-The-Wild Robot Teaching Without In-The-Wild Robots. In *Proceedings of Robotics: Science and Systems (RSS)*, 2024.
- [35] E. Todorov, T. Erez, and Y. Tassa. MuJoCo: A physics engine for model-based control. In *International Conference on Intelligent Robots and Systems (IROS)*, 2012.
- [36] Z. Fu, T. Zhao, and C. Finn. Mobile ALOHA: Learning Bimanual Mobile Manipulation with Low-Cost Whole-Body Teleoperation. *arXiv*, 2024.

Overview of Appendices

In the appendices below, we provide additional details on the implementation of RoboCrowd, our experiments, and our crowdsourced dataset. We provide a brief overview of each appendix below. For videos, please see our website: <https://robocrowd.github.io>

Appendix A – Task Details

We give descriptions of each of our 6 tasks, as well as renderings and images depicting sample expert demonstrations for each task.

Appendix B – Dataset Examples

We provide sample trajectories from our collected dataset including their task and quality annotations, to qualitatively illustrate the diversity of the behaviors in the dataset.

Appendix C – Additional Dataset Analysis

We provide further data analysis, including an offline user study to justify our scene choices, additional data quality analysis, and results on users’ self-reported Likert ratings of their interactions with the system.

Appendix D – Additional Details on Policy Learning Experiments

We provide additional details on the training and evaluation procedures for our policy learning experiments, as well as further qualitative analysis of the results.

Appendix E – Additional Details on Software Implementation and Data Annotation

We provide further details on the graphical user interface, interactive tutorial, software implementation, and data annotation pipeline.

Appendix F – Additional Details on Pilot Studies and System Development

We provide more details on how we designed and refined the system through pilot studies.

Appendix G – Overview of Action Chunking with Transformers (ACT) [1]

We provide additional background on the Action Chunking with Transformers (ACT) algorithm.

A Task Details

In [Tables 2 to 7](#) below, we provide a verbal description of the behavior that the expert demonstrations perform for each task. We additionally include a virtual rendering of different segments of a sample demonstration (where the gripper is rendered with increasing opacity for later timesteps). Additionally, we show a timelapse of the overhead camera image observation for the same sample expert demonstration.


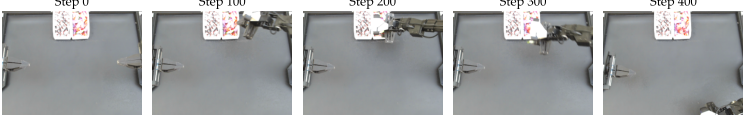

Task Name	Pick up a Hi-Chew (hi-chew)
Task Description	Move the right arm towards the candy bin. Grasp one Hi-Chew. Drop it in the End Zone. Finally, return to the home position.
Expert Trajectory Rendering	<div style="display: flex; justify-content: space-around; text-align: center;"> <div>Steps 0 → 249</div> <div>Steps 250 → 449</div> <div>Steps 450 → 504</div> </div> 
Expert Trajectory Timelapse	<div style="display: flex; justify-content: space-around; text-align: center;"> <div>Step 0</div> <div>Step 100</div> <div>Step 200</div> <div>Step 300</div> <div>Step 400</div> </div>  <div style="display: flex; justify-content: center; text-align: center; margin-top: 10px;"> <div>Step 500</div> </div> 

Table 2: Description of the hi-chew task, as well as a rendering and timelapse of a sample expert trajectory.


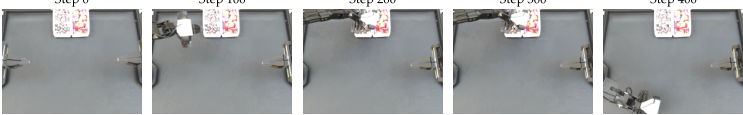
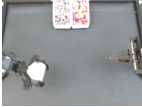
Task Name	Pick up a Tootsie Roll (tootsie-roll)
Task Description	Move the left arm towards the candy bin. Grasp one Tootsie Roll. Drop it in the End Zone. Finally, return to the home position.
Expert Trajectory Rendering	<div style="display: flex; justify-content: space-around; text-align: center;"> <div>Steps 0 → 249</div> <div>Steps 250 → 499</div> <div>Steps 500 → 599</div> </div> 
Expert Trajectory Timelapse	<div style="display: flex; justify-content: space-around; text-align: center;"> <div>Step 0</div> <div>Step 100</div> <div>Step 200</div> <div>Step 300</div> <div>Step 400</div> </div>  <div style="display: flex; justify-content: center; text-align: center; margin-top: 10px;"> <div>Step 500</div> </div> 

Table 3: Description of the tootsie-roll task, as well as a rendering and timelapse of a sample expert trajectory.

Task Name	Pick up a Hershey Kiss (hershey-kiss)
Task Description	Move the right arm or the left arm towards the candy bin. Grasp one Hershey Kiss. Drop it in the End Zone. Finally, return to the home position.
Expert Trajectory Rendering	<div style="display: flex; justify-content: space-around;"> <div>Steps 0 → 249</div> <div>Steps 250 → 399</div> <div>Steps 400 → 453</div> </div>
Expert Trajectory Timelapse	

Table 4: Description of the hershey-kiss task, as well as a rendering and timelapse of a sample expert trajectory.

Task Name	Eject a Jelly Bean from the Candy Dispenser (jelly-bean)
Task Description	Use the left arm to pull a cup from the cup dispenser. Bring the cup near the lever of the candy dispenser. Use the right arm to align the cup under the lever, then press the lever. Then, use the right arm to pick up the cup and bring it to the End Zone. Finally, return to the home position.
Expert Trajectory Rendering	<div style="display: flex; justify-content: space-around;"> <div>Steps 0 → 249</div> <div>Steps 250 → 499</div> <div>Steps 500 → 624</div> <div>Steps 625 → 874</div> <div>Steps 875 → 1049</div> </div> <div style="margin-top: 10px;"> <div>Steps 1050 → 1131</div> </div>
Expert Trajectory Timelapse	

Table 5: Description of the jelly-bean task, as well as a rendering and timelapse of a sample expert trajectory.

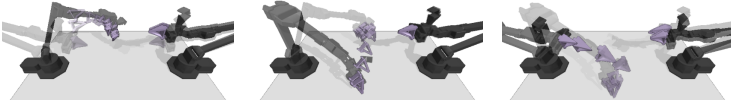


Task Name	Pick up a Hi-Chew from the Bin (hi-chew-bin)
Task Description	Move the right arm or the left arm towards the candy bin. Grasp one Hi-Chew. Drop it in the End Zone. Finally, return to the home position.
Expert Trajectory Rendering	<div style="display: flex; justify-content: space-around; text-align: center;"> <div data-bbox="602 835 732 858">Steps 0 → 249</div> <div data-bbox="841 835 987 858">Steps 250 → 549</div> <div data-bbox="1089 835 1235 858">Steps 550 → 741</div> </div> 
Expert Trajectory Timelapse	<div style="display: flex; justify-content: space-around; text-align: center;"> <div data-bbox="597 1012 638 1035">Step 0</div> <div data-bbox="743 1012 784 1035">Step 100</div> <div data-bbox="889 1012 930 1035">Step 200</div> <div data-bbox="1036 1012 1076 1035">Step 300</div> <div data-bbox="1182 1012 1222 1035">Step 400</div> </div>  <div style="display: flex; justify-content: space-around; text-align: center;"> <div data-bbox="597 1150 638 1173">Step 500</div> <div data-bbox="743 1150 784 1173">Step 600</div> <div data-bbox="889 1150 930 1173">Step 700</div> </div> 

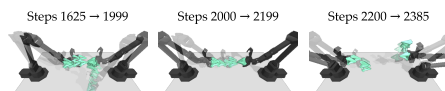
Table 6: Description of the hi-chew-bin task, as well as a rendering and timelapse of a sample expert trajectory.

Task Name Open the Ziploc, Pick up a Hi-Chew, then Close the Ziploc (hi-chew-ziploc)

Task Description Use the right arm to bring the Ziploc bag to the center of the table. Then, use the left arm to hold the Ziploc while pulling the Ziploc tab with the right arm to open the bag. Then, spread the Ziploc open and pick out a Hi-Chew with the right arm, and bring it to the End Zone. Then, use the right arm to hold the Ziploc while pulling the Ziploc tab closed with the left arm. Finally, use the right arm to place the Ziploc back in the corner of the table, and return the arms to the home position.



Expert Trajectory
Rendering



Expert Trajectory
Timelapse

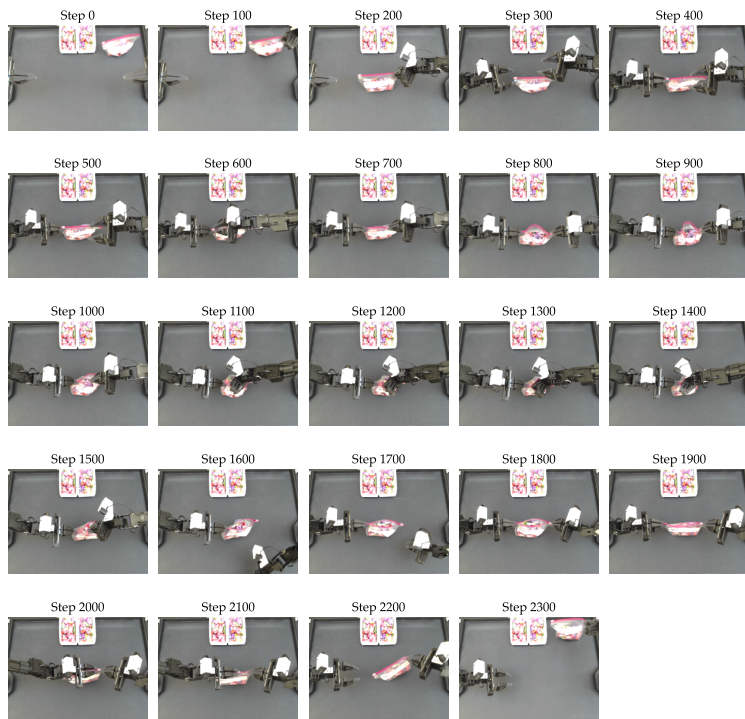


Table 7: Description of the hi-chew-ziploc task, as well as a rendering and timelapse of a sample expert trajectory.

B Dataset Examples

In [Figs. 5 to 7](#), we give 3 qualitative examples of interaction episodes in our crowdsourced dataset. We illustrate a timelapse of each episode with the overhead camera observation. We also include the task and quality annotations at each timestep, with a verbal description of the episode in the caption.

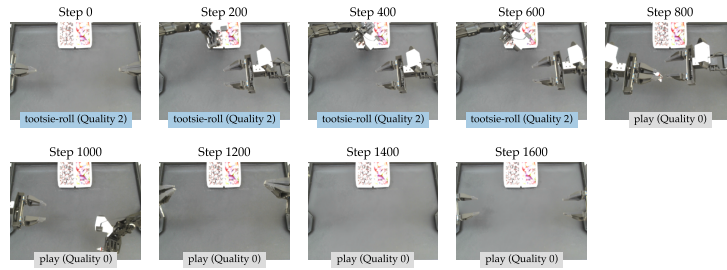


Figure 5: In this trajectory, the user begins by performing the `tootsie-roll` task with moderate quality—i.e., there are about 3 attempts to grasp the candy, and there is some extraneous movement in the right arm, but the user is otherwise successful at grasping the candy. Before bringing the candy all the way to the End Zone, the user attempts to unwrap the candy. They then hand it over to the other arm, place it in the End Zone, and then move the arms upward. The first half of the episode is marked as `tootsie-roll` (Quality 2) and the latter half of the episode is marked as `play` (Quality 0).

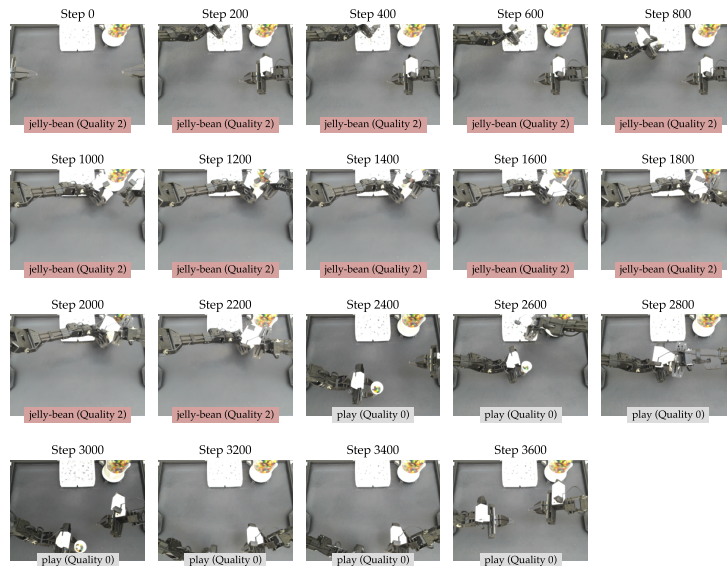


Figure 6: In this trajectory, the user grasps a cup from the cup dispenser and places it under the lever of the candy machine. They are successful in collecting jelly beans in the cup, though the trajectory includes retrying behavior and is not as smooth as an expert trajectory. The user brings the cup halfway to the End Zone, and then begins behaviors that are not part of the task—i.e., placing a Hershey Kiss in the cup before bringing it to the End Zone. The first part of the episode is marked as `jelly-bean` (Quality 2) and the latter part is marked as `play` (Quality 0).



Figure 7: In this trajectory, the user correctly moves the Ziploc from the corner of the table to the center of the table, and grasps a Hi-Chew from inside the Ziploc which they bring to the End Zone. They are unsuccessful in closing the Ziploc before episode termination. The user is task-directed for the whole episode, however takes longer than better quality trajectories for this task and performs retrying behavior at each subtask. The whole trajectory is marked as hi-chew-ziploc (Quality 1).

C Additional Dataset Analysis

In this section, we provide additional data analysis. In [Appendix C.1](#), we describe an offline study over user preferences for different candies, informing our different scene setups. In [Appendix C.2](#) and [Appendix C.2.1](#), we examine additional metrics (i.e., tutorial quality and Likert ratings) that correlate with quality of user interaction episodes, and in [Appendix C.3](#), we provide additional statistics on usage and retention.

C.1 Justification for Scene Choices

To justify our scene setup and task pairings, we perform an offline survey on user preferences for various candies. On a sample of $N = 16$ users, we find that 81% prefer a Hi-Chew to a Tootsie Roll. Thus, BinScene (which includes the hi-chew and tootsie-roll tasks) allows us investigate whether this preference for material reward shapes task choice when teleoperating demonstrations, when the task is otherwise equivalent besides the material reward. Users exhibit a more mild preference for a Hershey Kiss compared to a small handful of Jelly Beans (with 62% of respondents preferring the Hershey Kiss). Bin+ZiplocScene (which includes the hi-chew-bin and hi-chew-ziplloc tasks) allows us to investigate how intrinsic motivation and task difficulty affects user behavior when teleoperating in the case that the material reward (a Hi-Chew) is held constant between the the simpler task and the more challenging task. Bin+DispenserScene allows us to investigate this question when the material rewards are different, and users do not exhibit an overall preference for the reward from the harder task (and even mildly prefer the reward from the easier task).

C.2 Additional Metrics on Demonstration Quality

Our crowdsourced dataset contains rich interaction data per user ID—during and after the interactive tutorial period. This dataset can help to yield insights about which users give higher quality trajectories, and what factors can help predict this quality. As an example, we examine how the quality of interactions *after* the tutorial (i.e., when the user selects tasks in the scene to perform) correlates with quality *during* the tutorial period (i.e., when the user is instructed to complete simple onboarding tasks). Specifically, we examine the distribution of mean quality during task interactions versus minimum quality during the tutorial period; the user’s tutorial period is classified as 0 if there is any off-task behavior, 1 if the tutorial is performed but with retrying, and 2 if the tutorial is performed smoothly. We observe a loose positive correlation between higher minimum tutorial quality and mean task quality; and notably, users who produce consistently high quality task demonstrations (quality 3) are more present in the group with high quality tutorials. The tutorial period can therefore be a first-cut proxy at filtering demonstrators by quality.

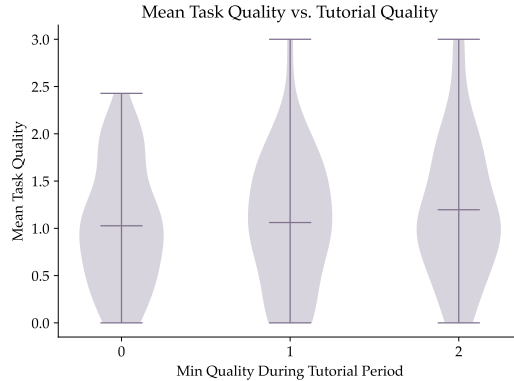


Figure 8: Distribution of Mean Task Quality versus Minimum Quality during the Tutorial Period.

C.2.1 Self-Reported Likert Metrics

After every interaction episode, we prompt the user to answer whether they agree with 3 statements, on a 5-point scale (1 - Strongly Disagree; 2 - Disagree; 3 - Neutral; 4 - Agree; 5 - Strongly Agree).

- Intuitive: Controlling the robot was intuitive.
- Interesting: Controlling the robot was fun and interesting.
- Wanted: The robot accomplished the task in the way that I wanted.

Fig. 9 summarizes the responses to these questions, aggregated by users’ minimum ratings to each statement over their interaction episodes. The majority of users agree with all three statements, and most often have the strongest ratings for Interesting compared to Intuitive and Wanted. We find also that there are loose correlations between the manually annotated quality scores for users’ interaction episodes and users’ self-reported ratings for each of these metrics. Specifically, users who self-report low ratings on each of the three metrics have lower mean quality scores. However, users who self-report high ratings have quality scores that span low to high.

C.3 Usage and Retention

We illustrate the usage of the RoboCrowd in Fig. 10. We observe significant engagement with RoboCrowd over the two-week collection period: there were $N = 231$ unique users in total. On most days, more than two-thirds of these were new users that had not used the system on prior days. There were a total of 817 interaction episodes distributed throughout the period. The most common time at which users interacted with the system was about 1pm, corresponding to the most trafficked time in the café (lunchtime). We collect 129 interaction episodes in BinScene (Day 1), 381 in Bin+DispenserScene (Days 2-5), and 307 in Bin+ZiplocScene (Days 6-11).

D Additional Details on Policy Learning Experiments

In this section, we give additional details on our policy learning experiments. Appendix D.1 provides training details and hyperparameters, Appendix D.2 provides details on our evaluation procedure, and Appendix D.3 provides additional qualitative discussion of our learned policies.

D.1 Training Details

For the *Expert* and *Co-train* experiments, we train policies for 200K steps for all tasks. For the *Fine-tune* experiments, we fine-tune the co-trained model (partially trained for 100K steps) for an additional 50K steps on expert data only. We use the implementation of ACT [1] from [36], including the default hyperparameters from [1], as shown in Table 8.

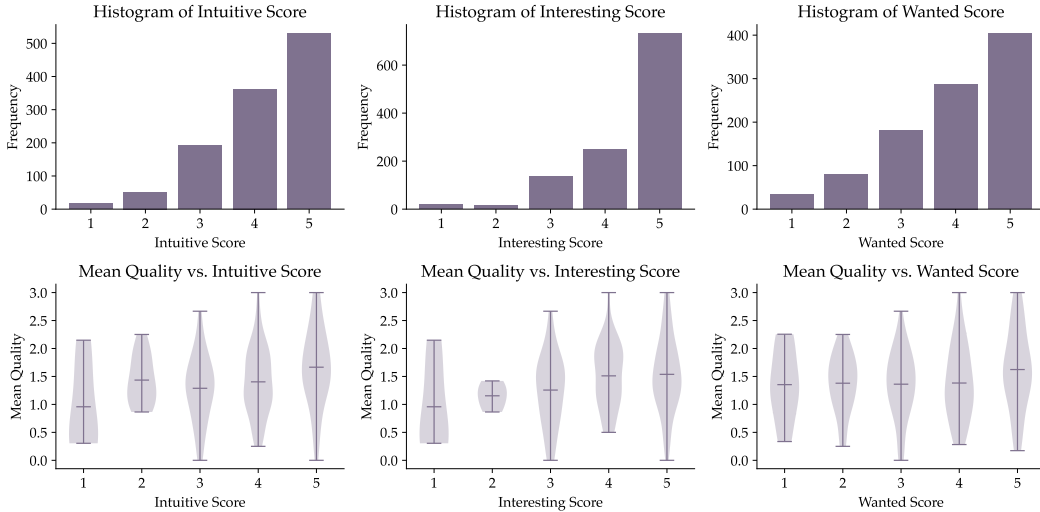


Figure 9: (*Top*) Histogram of Likert Ratings (aggregated by the user’s minimum response over their interaction episodes) for the Intuitive, Interesting, and Wanted questions. (*Bottom*) Distribution of mean quality of interaction episodes for different Likert Ratings for Intuitive, Interesting, and Wanted.

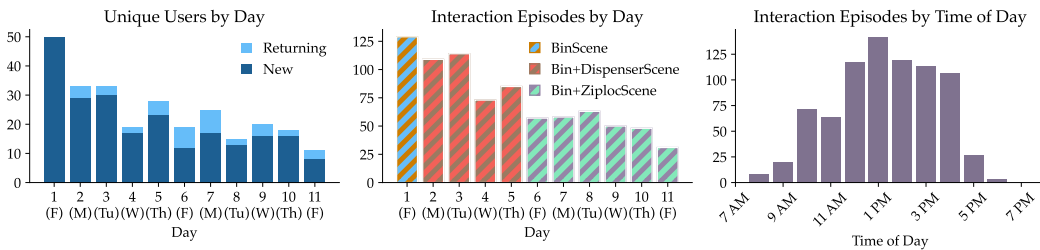


Figure 10: Statistics on usage over a two-week period: number of users per day (left), number of interaction episodes per day (middle), and distribution of interaction episodes by time of day (right).

D.2 Evaluation Details

We perform policy evaluations for 40 trials each, early stopping when policies exhibit excessively jittery or unsafe behavior. While the RoboCrowd training dataset was collected in a café where lighting varies throughout the day, during evaluation, we move the setup to a location with a visually similar background but consistent lighting for controlled evaluations.

For the bin-picking tasks, we define success as the robot arm picking exactly one of the desired candy and bringing it to the End Zone. For our challenging, long-horizon tasks (jelly-bean and hi-chew-ziploc), success is 0% for all policies, so we instead compare policies via normalized return to measure partial proficiency at tasks. We describe the process for computing normalized return below.

Each of the following subtasks in `jelly-bean` corresponds to 1 point in the episode return: Retrieves Cup from Dispenser; Places Cup Down; Aligns Cup Under Lever; Presses Lever; Collects Jelly Beans in Cup; Picks up Cup; Brings Cup to End Zone. Each of the following subtasks in `hi-chew-ziploc` corresponds to 1 point in the episode return: Picks up Bag; Places Bag in Center of Table; Slides Open; Picks Hi-Chew; Brings Hi-Chew to End Zone; Closes Bag; Places Bag in Corner of Table. For these tasks, we report normalized return—the average return over evaluation trials divided by the maximum return (achieved by all expert demonstrations).

Learning Rate	1e-5
Batch Size	8
# Encoder Layers	4
# Decoder Layers	7
Feedforward Dimension	3200
Hidden Dimension	512
# Heads	8
Chunk Size	100
KL-weight (β)	10
Dropout	0.1
Backbone	ResNet-18
Image Augmentations	RandomCrop, Random- Resize, RandomRotation, ColorJitter

Table 8: Hyperparameters for ACT, shared for all experiments.

D.3 Qualitative Analysis of Learned Policies

We find that in most cases, *Co-train* and/or *Fine-tune* improve upon *Expert*. However, the specific effects vary by task. For example, we find that for the *hi-chew* task, the co-trained policy performs worse than the expert policy, but the fine-tuned policy performs better; whereas with the *hershey-kiss* task, both the co-trained policy and fine-tuned policy perform better. We hypothesize that the crowdsourced data is more useful for *hershey-kiss* because (a) *hershey-kiss* is a more complex task (in that it is more multimodal, i.e., either arm can be used to pick up a Hershey Kiss, and the grasping required needs to be more precise to not crush the Hershey Kiss) and (b) a greater proportion of the *hershey-kiss* data is of higher quality. We notice that the crowdsourced data for *jelly-bean* is especially diverse, and naïvely co-training or fine-tuning underperforms using the expert data only.

Qualitatively, we observe in several cases that the co-trained and fine-tune policies exhibit meaningful but suboptimal behaviors from the crowdsourced data (e.g., picking up multiple objects from the bin instead of one). On the other hand, there are also helpful behaviors from the crowdsourced data (*not* represented in the expert data) that benefit trained policies—e.g., regrasping behavior.

Overall, the RoboCrowd dataset is very diverse, and contains both task-relevant behaviors (of various levels of quality) and free-play behavior. Future work on more sophisticated policy learning methods that leverage these diverse characteristics can help to get the maximum utility out of crowdsourced demonstration data.

E Additional Details on Software Implementation and Data Annotation

In this section, we provide additional details on our software interface and implementation, as well as our data annotation pipeline. [Appendix E.1](#) provides an overview of the application flow and interface, [Appendix E.2](#) details the interactive tutorial procedure, [Appendix E.3](#) provides implementation details, and [Appendix E.4](#) details the data annotation pipeline.

E.1 Application Flow and User Interface

[Fig. 11](#) gives an overview of the flow through the tablet application, and [Table 9](#) provides screenshots of the major pages referenced in the flowchart. We additionally highlight the Interactive Tutorial in [Fig. 12](#) and the visual warning for collision detection in [Fig. 13](#). We now briefly describe the application flow. To begin a new session, the user taps their ID card on the card reader, which advances the tablet application to a screen where the user can enter a nickname (if they are a new user). They are then directed to the Main Page, where they complete a consent form and the interactive tutorial. From the Main Page, users can also press a “Start Playing” button which directs them to the Task Page, where they can see videos of tasks available in the scene, and can tap on a task to see more details and begin demonstrating the task.

For safety, the user receives an audial and visual warning (Fig. 13) if the arms are near-collision. When users are done with the task (i.e., they click a Stop button on the Task Detail Page or they rest the grippers on the mechanical stop), they are asked to mark their demonstration as a success or failure, and fill out a brief survey. The success/failure markings are used as the basis for the points which are added to the user’s point total in the Leaderboard, which is accessible from the Main Page; in our experiments, users receive 10 points for successful “easy” tasks (bin-picking) and 20 points for successful “difficult” tasks (the remaining tasks). From the Main Page, users can also choose to provide feedback, or press a Request Help button which immediately notifies the study team (e.g., if the user needs assistance or if the setup requires maintenance).

E.2 Interactive Tutorial

We provide a zoomed-in version of the pages in the Interactive Tutorial in Fig. 12. The aim of the tutorial is to guide the user on how to start and stop interaction episodes as well as how to puppeteer with ALOHA. Specifically, users are first instructed to wait until ALOHA’s arms rise to the home position, and then they are given instructions on how to start puppeteering (by squeezing both sets of grippers on the leader arms). After they do so, the tutorial automatically proceeds to the next stage, where users then are told to gently touch the left and right arms to the table; the goal is to help users get calibrated to the robot’s range of motion and degrees of freedom, as well as the types of forces they need to apply to move the arms. Finally, users are given instructions on how to stop the interaction episode, by resting the grippers of the leader arms in the grooves of the mechanical stops. When the user does so, the puppet arms are automatically lowered, and the user is presented a brief video on how to navigate the rest of the interface.

E.3 Implementation Details

The software application is implemented with React (frontend) and Flask (backend), and uses WebSocket connections to communicate between the user client and backend server. We use a SocketIO-ROS bridge to pass messages between the backend server and robot controller. The robot controller operates at 50Hz and is based on [1]. When the robot is being teleoperated, we run a parallel simulation in MuJoCo [35] which is updated at every time step to detect self-collisions.

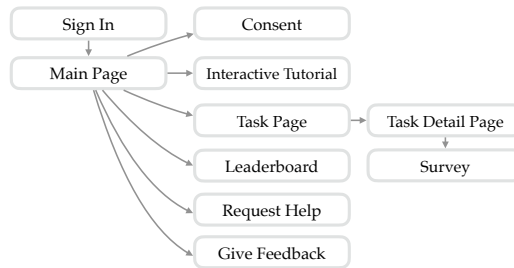


Figure 11: Flowchart illustration of pages in the user interface.

E.4 Data Annotation Pipeline

We annotate episodes in our crowdsourced dataset by task and quality. We implement an interface for annotation, which we illustrate in Fig. 14. We annotate episodes by dragging a slider which scrubs through the episode and selecting a task and quality annotation for different segments of the episode. We describe the annotation rules below.

- play (Quality 0). All free-play behavior is marked as play with quality 0. Play data includes undirected movements and tasks that the user makes up (e.g., trying to unwrap a candy). It also includes extraneous movements before and after the user performs a task.

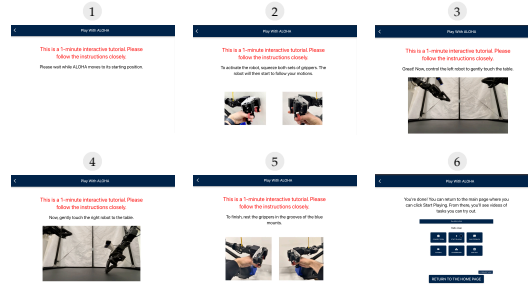


Figure 12: Screenshot of the pages in the interactive tutorial interface.

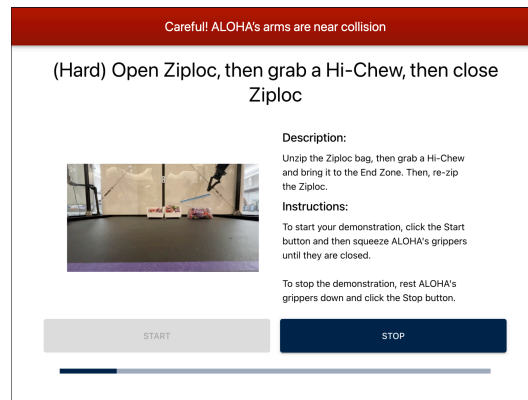


Figure 13: Screenshot of a visual collision warning on the task page. An auidial alarm (beeping sound) is played on the tablet when the visual collision warning appears.

- **tutorial** (Quality 1–2). Movements associated with the tutorial (e.g., touching the grippers to the table) are marked as Q1 if there is any retrying behavior and Q2 if the motions are smooth.
- **<task>** (Quality 1–3). Task-relevant motions for each of our six tasks are labeled with the task name and a quality from 1 to 3. Q3 is used to describe segments that complete subtasks smoothly with no more than 2 retries. Q2 is used to describe segments that use no more than 4 retries for any one subtask, or that are completed but with slight errors (e.g., grabbing more than 1 candy from a bin). Q1 is used to describe segments that are task-relevant but of poor quality (e.g., more than 4 retries for any one subtask), cause changes to the scene (e.g., dropping a candy on the table), or complete the task in a significantly different manner than the expert demonstrations (e.g., using the opposite arm for any subtask).

F Additional Details on Pilot Studies and System Development

Prior to full system deployment, we conducted pilot studies on a smaller population to help us iterate on our system. We obtained the Institutional Review Board’s approval before both the pilot studies and the full deployment. We recruited $N=10$ participants to interact with the system. In order to mimic organic interactions as closely as possible, we did not provide the participants with any verbal instructions, other than to begin interacting with the system as if they happened upon it organically. Our software interface guided the participants through the consent form and tutorial. Here is a sample of feedback provided by participants, coupled with changes we made to the system.

- *Degrees of Freedom*: Users indicated that puppeteering demonstrations was challenging the first time because they needed to “understand the degrees of freedom” of the robot. To address this feedback,

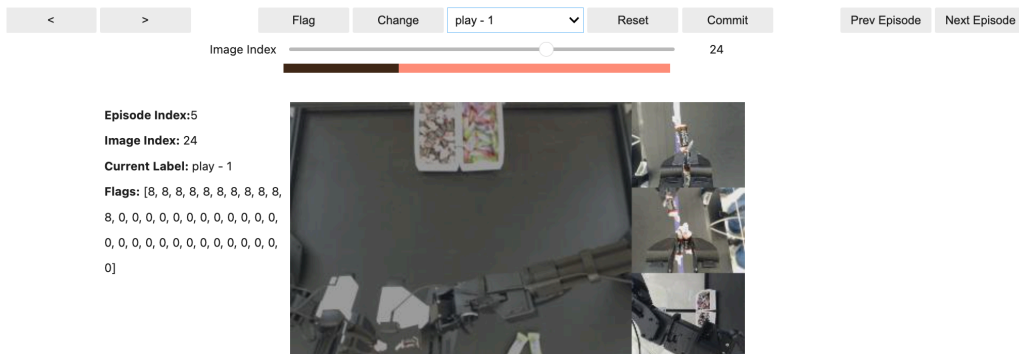


Figure 14: Screenshot of the data annotation interface. Annotators can scrub through the episode and label segments with task and quality labels, which color codes a bar to visualize the different tasks and qualities in the episode. When the annotator is done labeling an episode, they can “commit” their labels and proceed to the next episode.

we created a tutorial where the user was guided through how to perform primitive movements of the leader arms (e.g., controlling both puppet arms to touch the bottom of the workspace) before they began interacting with the system.

- *Tutorial Format*: In an initial prototype, our tutorial was a video that a user would watch before using the system. Users provided feedback that they felt “impatient” and would rather “explore what it is like to interface with the robot” rather than “watch a long video.” To address this feedback, we made the tutorial efficient and interactive: 4 steps that the user would perform with the robot after watching them on the screen. The interactive tutorial automatically advances after detecting that each step is complete.
- *Start and Stopping Demonstrations*: In an initial prototype, users begin demonstrations by (1) tapping a Start button on an interface and (2) squeezing the grippers of the leader arms closed. To terminate episodes, they would simply need to (1) leave the arms to rest on the robot body and (2) tap a Stop button on the interface. We received feedback that squeezing the gripper to start episodes “made sense” but the “rest position at the end was confusing.” To address this feedback, we designed and 3D printed a mechanical stop for users to rest the arms. We automatically terminate episodes when handles of the leader arms make contact with this mechanical stop.
- *Interface*: In an initial prototype, users would access the interface on their own smartphone by scanning a QR code pasted on the platform. A user reported that they would prefer if more of their interaction would happen “in the position that they will be doing the task.” We therefore switched to a tablet interface mounted at the base of the platform, which was accessible when the user sat down to begin interacting with the robot. On the interface itself, users reported that it was “easy to understand.”
- *Collisions*: We observed that participants did not actively pay much attention to collisions between the robots, as well as the collision of wrist-camera mounts and objects mounted on the table. To address this, we (1) added collision avoidance between the arms and the table, (2) added an audio-visual alarm when arms were near collision, and (3) mounted objects to the table so that they would not move.

G Overview of Action Chunking with Transformers (ACT)

In this section, we provide a more extended background overview of imitation learning (IL) and the Action Chunking with Transformers (ACT) algorithm [1].

Imitation learning (IL) aims to learn a policy π_θ parameterized by θ given access to a dataset \mathcal{D} composed of expert demonstrations. Defined within the framework of a standard partially observable Markov

decision process (POMDP), each trajectory $\xi \in \mathcal{D}$ is a sequence of observation-action transitions $\{(o_0, a_0), \dots, (o_T, a_T)\}$. Most commonly, IL is instantiated as behavior cloning, which trains π_θ to minimize the negative log-likelihood of data, $\mathcal{L}(\theta) = -\mathbb{E}_{(o,a) \sim \mathcal{D}} [\log \pi_\theta(a|o)]$.

In practice, the human-collected demonstrations in \mathcal{D} may be diverse. To effectively learn from such diverse data, we can condition the policy on a latent variable z , which helps to capture the variability in the demonstrations by representing different modes of behavior. Representing this policy as the decoder in a conditional variational autoencoder (cVAE), we in addition learn an encoder q_ϕ from (observation, action) pairs to the latent space: $q_\phi(z | a^t, o^t)$. And we condition our policy on the latent variable: $\pi_\theta(\hat{a}^t | o_t, z)$. At test time, we sample latent vectors from the standard normal distribution, $z \sim \mathcal{N}(0, 1)$. We regularize the outputs of our encoder towards this distribution via a KL-penalty: $D_{KL}(q_\phi(z | a^t, o^t) \| \mathcal{N}(0, 1))$. This method is formalized as Action Chunking with Transformers (ACT) [1], an imitation learning algorithm designed to learn from diverse human demonstrations.




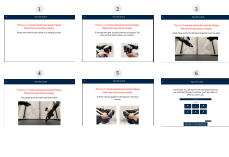
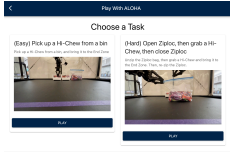
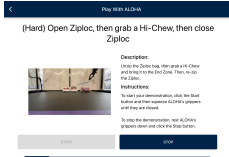
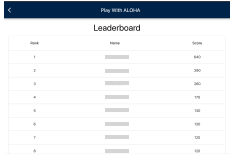
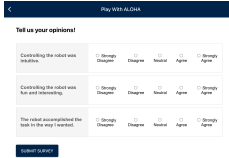


Page Name	Screenshot	Page Name	Screenshot
Sign In (Tap ID Card)		Sign In (Create User Profile)	
Main Page		Interactive Tutorial	
Task Page		Task Detail Page	
Leaderboard		Survey Page	
Request Help		Give Feedback	

Table 9: Screenshots of pages in the user interface.