

OPEN TEACH: A Versatile Teleoperation System for Robotic Manipulation

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Open-sourced, user-friendly tools form the bedrock of scientific ad-
2 advancement across disciplines. The widespread adoption of data-driven learning
3 has led to remarkable progress in multi-fingered dexterity, bimanual manipulation,
4 and applications ranging from logistics to home robotics. However, existing data
5 collection platforms are often proprietary, costly, or tailored to specific robotic
6 morphologies. We present OPEN TEACH, a new teleoperation system leveraging
7 VR headsets to immerse users in mixed reality for intuitive robot control. Built on
8 the affordable Meta Quest 3, which costs \$500, OPEN TEACH enables real-time
9 control of various robots, including multi-fingered hands, bimanual arms, and mo-
10 bile manipulators, through an easy-to-use app. Using natural hand gestures and
11 movements, users can manipulate robots at up to 90Hz with smooth visual feed-
12 back and interface widgets offering closeup environment views.

13 **Keywords:** Low Cost Teleoperation, Data Collection, Imitation learning

14 1 Introduction

15 The integration of learning-based methods has sparked a revolution in robotics, significantly enhanc-
16 ing capabilities in manipulation [1, 2, 3, 4], locomotion [5, 6, 7, 8], and aerial robotics [9, 10, 11].
17 More recent work has been making advancements in complex single-task behavior learning [12,
18 13, 2], multitask scenarios [14, 15], multimodal behavior learning [16, 17, 18, 19], and efficient
19 fine-tuning of pretrained behavior models [20, 21, 22]. A fundamental requirement across all these
20 threads of research is the need to collect data in the form of task demonstrations.

21 Commonly used teleoperation systems include devices such as joysticks and 3D spacemouses [23,
22 24], commercial VR headsets [25, 26, 27, 13, 28, 29], kinesthetic teaching [30], and phone tele-
23 operation [31]. Recently proposed exoskeleton-based teleoperation frameworks like ALOHA [2],
24 GELLO [32], and AirExo [33] attempt to alleviate this problem by having the human teleopera-
25 tor directly control a kinematically isomorphic version of the same robot arm. These frameworks
26 directly impose the kinematic constraints of the robot arm during teleoperation making it more com-
27 patible and intuitive to control the robot. Although highly effective, these systems can require an
28 additional robot for each robot being controlled, have high initial setup costs, and are designed for
29 specific robot morphologies. The challenge of easy-to-use teleoperation devices is more apparent
30 in dexterous manipulation problems [34, 35, 27, 13], owing to the high dimensional action space.
31 Such frameworks typically involve the use of expensive gloves [36, 37, 38], extensive calibration
32 processes [34, 27], or are susceptible to monocular occlusions [27].

33 In this work, we present OPEN TEACH, an open-source framework for robot teleoperation that sup-
34 ports a variety of robots, including bimanual and multi-finger manipulation, all at a price of \$500.
35 OPEN TEACH uses a VR headset (e.g. Quest 3) to put users / teachers in an immersive virtual world
36 where they can view a robotic scene both through their eyes, via visual passthrough, as well as re-
37 altime streams from the robot’s cameras. To control the robot, users can simply use hand gestures,
38 which are detected using onboard hand-pose estimators at 90Hz. We experimentally evaluate OPEN
39 TEACH on 38 tasks across single arm, bimanual, multi-fingered, and mobile manipulation robot se-
40 tups in both simulation and the real world. The tasks span from tabletop manipulation to contact-rich

41 dexterous manipulation. cross different robot morphologies, we find that users can provide demon-
42 strations at speeds on par with robot-specific teleoperation systems and significantly faster than
43 general-purpose systems like AnyTeleop [35]. Importantly, policies trained on the data collected
44 achieve an average success rate of 86% on 10 tasks in simulation and the real world, validating
45 the utility of policy learning using OPEN TEACH. The contributions of this work is summarized as
46 follows:

- 47 1. We present OPEN TEACH, an open-source system for plug-and-play teleoperation frame-
48 work suitable for collecting demonstrations across different robot morphologies in both
49 simulation and the real world.
- 50 2. We experimentally show that the demonstrations collected by OPEN TEACH can be used to
51 train effective, general-purpose manipulation behaviors.
- 52 3. Our user study on 15 new users highlights the efficacy of OPEN TEACH for both experi-
53 enced and new users.

54 OPEN TEACH will be fully open-sourced with mixed reality API, policy training code, and demon-
55 strations collected using OPEN TEACH available at <https://anon-open-teach.github.io/>.

56 2 OPEN TEACH

57 In OPEN TEACH, a user wears a Virtual Reality (VR) headset to provide demonstrations to a robot.
58 This involves creating a virtual world for teaching, retargeting the teacher’s hand and wrist pose to
59 the robot joints, and finally controlling the robot. We compare OPEN TEACH with various other
60 teleoperation systems across a variety of robot types and observe that OPEN TEACH is the only
61 framework that enables controlling multiple arms, hands, and mobile manipulators, is calibration-
62 free, and is completely open-source.

63 2.1 Placing an Operator in a Virtual World

64 We use the Meta Quest 3 VR headset to place the human teacher in a virtual world. The headset
65 surrounds the human in a virtual environment at a resolution of 2064×2208 and a native refresh
66 rate of 90Hz. The base version of this headset is affordable at \$499 and is relatively light at 513g.
67 Compared to the Meta Quest 2 VR headset used in prior work [13], the Quest 3 provides a full-color
68 passthrough allowing the human to get a direct view of the robot setup during teleoperation. These
69 features, especially the full-color passthrough, allow for a comfortable and intuitive operation by the
70 user. Additionally, similar to Arunachalam et al. [13], the Quest 3 API interface allows for creating
71 custom mixed reality worlds that visualize the robotic system along with diagnostic panels in VR. It
72 is important to highlight the exceptional clarity of the scene passthrough visible in Quest 3.

73 2.2 Pose Estimation with VR Headsets

74 Similar to Arunachalam et al. [13], we directly use the in-built hand pose estimator [39] of the
75 Quest 3 using 2 monochrome cameras. This is significantly more robust compared to single camera
76 alternatives [40]. Further, since the cameras are internally calibrated, they do not require specialized
77 calibration routines that are needed in prior multi-camera teleoperation frameworks [34, 35]. Also,
78 since the hand-pose estimator is integrated into the device, it can stream real-time hand poses at
79 90Hz.

80 2.3 Human to Robot Pose Retargeting

81 The inbuilt hand pose estimate from the VR headset provides us with the joint positions of all the
82 fingers of the human hand and the wrist. With this information, we can design wrappers that use
83 combinations of these joint positions to map the human hand poses to the robot poses for any given
84 robot morphology. In this work, we use a variety of robot arms, each with either a 2-fingered gripper
85 or a multi-fingered robot hand.

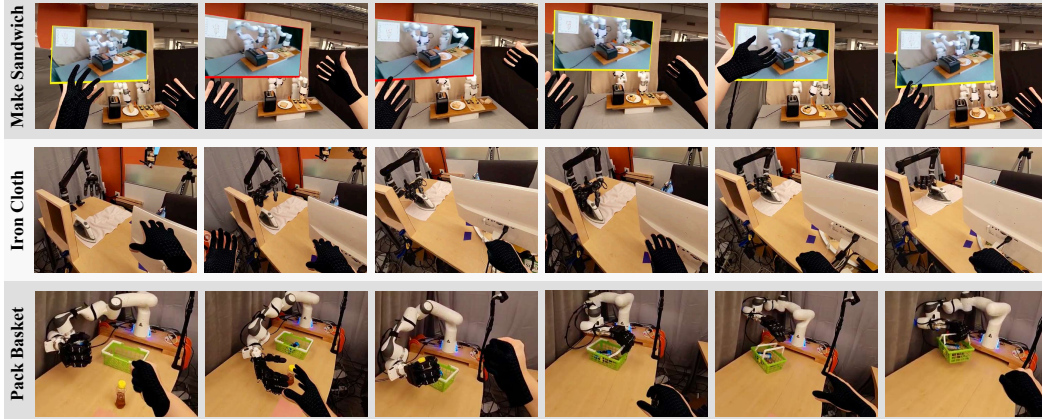


Figure 1: The demonstration collection process as viewed from within the VR application. Shown here is one task being performed for each real-world setup. High resolution images streamed at 90 Hz to the VR application allow for an immersive experience and enable reactive control by the user.

86 **Robot Arm:** We establish a 3D coordinate system using the wrist keypoint and knuckle points of the
 87 index and pinky fingers to define a 2D plane along the palm and a perpendicular third axis. The wrist
 88 position maps to the robot end effector position. Changes in the orientation of this hand coordinate
 89 system over time map to adjustments in end effector orientation.

90 **Robot Hand:** We use the teacher’s hand pose obtained from the VR to compute the individual joint
 91 angles in the teacher’s hand. Given these joint angles, a straightforward method of retargeting is
 92 to directly command the robot’s joints to the corresponding angles. In practice, this works well for
 93 all fingers except the thumb. To address this, we improve upon Arunachalam et al. [13], where the
 94 spatial coordinate of the teacher’s thumb tip is mapped to that of the robot hand and then an inverse
 95 kinematics solver is used to compute the joint angles of the thumb.

96 **Two-fingered gripper:** To detect the opening and closing of the two-fingered gripper, we utilize the
 97 pinch between the pinky finger and the thumb. We use a toggle mechanism for opening and closing
 98 the gripper where each pinch indicates toggling to the alternate state of the gripper.

99 **Mobile manipulator:** The same 3D coordinate system established for controlling robot arms is used
 100 for mapping the wrist’s movements to actions of the mobile robot. When the wrist moves forward,
 101 it extends the robot’s arm, enabling it to reach farther. Vertical wrist movements adjust the robot’s
 102 height, while lateral wrist movements cause the robot to move sideways by controlling its wheels.

103 3 Experiments

104 We demonstrate the usefulness of the collected data by training visual and visuotactile policies using
 105 behavior cloning [41] and inverse RL [42, 43].

106 3.1 Imitation Learning with OPEN TEACH Data

107 Here, we describe the algorithms used for learning policies on data collected through OPEN TEACH.

- 108 1. **Franka-Allegro:** We record both visual and tactile data for this setup. The policies are
 109 trained using TAVI [44], a demonstration-guided residual RL algorithm that collects a few
 110 expert demonstrations and learns a robot policy using both visual and tactile data.
- 111 2. **Allegro Sim:** We only record visual data for this setup and train policies using FISH [21].
- 112 3. **LIBERO Sim [23]:** We only record visual data for this setup. The policies are trained
 113 using transformer-based BC with a GMM head [45] and action chunking [2].

Table 2: Performance of policies learned on data collected through OPEN TEACH. For Franka-Allegro, Allegro Sim, and Libero Sim, TAVI [44], FISH [21] and BC were respectively used to train policies.

| Domain | Robot Setup | Stream Frequency (in Hz) | |
|--------|----------------|--------------------------|--------------|
| | | Arm | End Effector |
| Real | Franka-Allegro | 60 | 60 |
| | Kinova-Allegro | 60 | 60 |
| | Bimanual | 90 | 90 |
| | Stretch | 5 | 5 |
| Sim | Allegro Sim | 60 | 60 |
| | LIBERO Sim | 20 | 20 |

| Robot Setup | Task | Number of Demos | Success Rate |
|----------------|---------------------------------|-----------------|--------------|
| Franka-Allegro | Open Box | 3 | 9/10 |
| | Grasp Sponge | 6 | 7/10 |
| | Pick Up Tea Sachet | 4 | 7/10 |
| Allegro Sim | Grasp Object and Twist | 6 | 8/10 |
| | Flip Cube | 6 | 10/10 |
| | Flip Sponge | 6 | 10/10 |
| Libero Sim | Pinch Grasp | 6 | 7/10 |
| | Close Top Drawer of Cabinet | 10 | 10/10 |
| | Turn on Stove | 10 | 9/10 |
| | Pick and Place Soup into Basket | 50 | 9/10 |

Table 3: User study comparing OPEN TEACH with baselines when used by experts and new users.

| Task | Success Rate | | | | Median completion time for successful demonstrations (in s) | | | |
|-------------------|--------------|-----------|------------|------------|---|-----------|------------|------------|
| | New User | | | Expert | New User | | | Expert |
| | Holo-Dex | AnyTeleop | Open Teach | Open Teach | Holo-Dex | AnyTeleop | Open Teach | Open Teach |
| Flip cube | 1 | 1 | 1 | 1 | 6.58 | 13.71 | 5.5 | 2.85 |
| Pinch Grasp | 0 | 0.2 | 0.8 | 1 | 17.49 | 18.94 | 18.72 | 3.71 |
| Pour | N/A | N/A | 0.4 | 0.8 | N/A | N/A | 40.97 | 14.83 |
| Pick and Place | N/A | N/A | 0.8 | 0.8 | N/A | N/A | 23.57 | 11.875 |
| Open box of mints | N/A | N/A | 0.5 | 1 | N/A | N/A | 32.21 | 20.45 |

114 3.2 How versatile is OPEN TEACH across robotic setups?

115 The primary idea behind OPEN TEACH is that given any robotic setup, a user can purchase an
 116 affordable off-the-shelf VR headset (in this case, Quest 3) and plug the headset and robot setup into
 117 the proposed framework to start teleoperating the robot without any additional hardware setup cost.
 118 To investigate its versatility, we use OPEN TEACH to teleoperate four different real world robotic
 119 setups, each having a different combination of a robot arm and end effector type — Franka Allegro,
 120 Kinova Allegro, a Bimanual setup with 2 xArm7 robots, and Hello Stretch for mobile manipulation.
 121 We also exhibit the applicability of OPEN TEACH in simulation through evaluations on 2 simulated
 122 environment suites — Allegro Sim and LIBERO Sim [23]. The frequency of teleoperation for each
 123 of the setups has been provided in Table 1.

124 3.3 How successful are policies trained with OPEN TEACH?

125 Table 2 provides the success rates of policies learned using imitation learning across both the real-
 126 world and simulated setups. We use TAVI [44] to learn visuotactile policies on Franka-Allegro, and
 127 FISH [21] to learn visual policies on Allegro Sim. Similar to prior work [44, 21], these policies were
 128 learned within 20 minutes and achieved an average success rate of 82%, validating the high quality
 129 of the collected observation data. Behavior cloning policies on LIBERO Sim achieve an average
 130 success rate of 93%, confirming the high quality of the collected action data. Overall, the learned
 131 policies achieve an average success rate of 86% across all tasks and robot morphologies.

132 4 Conclusion

133 In this work, we introduce OPEN TEACH, an open-source unified framework designed to facilitate
 134 low-latency, high-frequency robot teleoperation. This versatile framework is tailored to accommo-
 135 date diverse tasks and is compatible with a range of robot morphologies. However, we recognize
 136 a few limitations in this work: (a) OPEN TEACH relies on the accuracy of the in-built hand pose
 137 detection in the VR headset. Inaccuracies, particularly when fingers are occluded from view, can
 138 diminish the quality of hand tracking, posing challenges to teleoperation. (b) In specific instances,
 139 the pose detector on the Oculus board may misconstrue finger positions, leading to difficulties in
 140 executing gestures like gripper closing, which relies on precise pinches between fingers. Addressing
 141 these challenges through future research on hand pose detection and tracking holds the potential to
 142 enhance the ease and intuitiveness of teleoperation using VR headsets.

References

- 143
- 144 [1] C. Chi, B. Burchfiel, E. Cousineau, S. Feng, and S. Song. Iterative residual policy: for goal-
145 conditioned dynamic manipulation of deformable objects. *arXiv preprint arXiv:2203.00663*,
146 2022.
- 147 [2] T. Z. Zhao, V. Kumar, S. Levine, and C. Finn. Learning fine-grained bimanual manipulation
148 with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.
- 149 [3] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Haus-
150 man, A. Herzog, J. Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv*
151 *preprint arXiv:2212.06817*, 2022.
- 152 [4] N. M. M. Shafiullah, A. Rai, H. Etukuru, Y. Liu, I. Misra, S. Chintala, and L. Pinto. On
153 bringing robots home. *arXiv preprint arXiv:2311.16098*, 2023.
- 154 [5] S. Gangapurwala, M. Geisert, R. Orsolino, M. Fallon, and I. Havoutis. Rloc: Terrain-aware
155 legged locomotion using reinforcement learning and optimal control. *IEEE Transactions on*
156 *Robotics*, 2022.
- 157 [6] Y. Ma, F. Farshidian, T. Miki, J. Lee, and M. Hutter. Combining learning-based locomotion
158 policy with model-based manipulation for legged mobile manipulators. *IEEE Robotics and*
159 *Automation Letters*, 7(2):2377–2384, 2022. doi:10.1109/LRA.2022.3143567.
- 160 [7] L. Smith, I. Kostrikov, and S. Levine. A walk in the park: Learning to walk in 20 minutes with
161 model-free reinforcement learning. *arXiv preprint arXiv:2208.07860*, 2022.
- 162 [8] X. Cheng, K. Shi, A. Agarwal, and D. Pathak. Extreme parkour with legged robots. *arXiv*
163 *preprint arXiv:2309.14341*, 2023.
- 164 [9] T. Zhang, G. Kahn, S. Levine, and P. Abbeel. Learning deep control policies for autonomous
165 aerial vehicles with mpc-guided policy search. In *2016 IEEE international conference on*
166 *robotics and automation (ICRA)*, pages 528–535. IEEE, 2016.
- 167 [10] D. Gandhi, L. Pinto, and A. Gupta. Learning to fly by crashing. In *2017 IEEE/RSJ Interna-*
168 *tional Conference on Intelligent Robots and Systems (IROS)*, pages 3948–3955. IEEE, 2017.
- 169 [11] J. Hwangbo, I. Sa, R. Siegwart, and M. Hutter. Control of a quadrotor with reinforcement
170 learning. *IEEE Robotics and Automation Letters*, 2(4):2096–2103, 2017. doi:10.1109/LRA.
171 2017.2720851.
- 172 [12] C. Wang, L. Fan, J. Sun, R. Zhang, L. Fei-Fei, D. Xu, Y. Zhu, and A. Anandkumar. Mimicplay:
173 Long-horizon imitation learning by watching human play. *arXiv preprint arXiv:2302.12422*,
174 2023.
- 175 [13] S. P. Arunachalam, I. Güzey, S. Chintala, and L. Pinto. Holo-dex: Teaching dexterity with
176 immersive mixed reality. In *2023 IEEE International Conference on Robotics and Automation*
177 *(ICRA)*, pages 5962–5969. IEEE, 2023.
- 178 [14] A. Padalkar, A. Pooley, A. Jain, A. Bewley, A. Herzog, A. Irpan, A. Khazatsky, A. Rai,
179 A. Singh, A. Brohan, et al. Open x-embodiment: Robotic learning datasets and rt-x mod-
180 els. *arXiv preprint arXiv:2310.08864*, 2023.
- 181 [15] H. Bharadhwaj, J. Vakil, M. Sharma, A. Gupta, S. Tulsiani, and V. Kumar. Roboagent: Gener-
182 alization and efficiency in robot manipulation via semantic augmentations and action chunking.
183 *arXiv preprint arXiv:2309.01918*, 2023.
- 184 [16] N. M. Shafiullah, Z. Cui, A. A. Altanzaya, and L. Pinto. Behavior transformers: Cloning k
185 modes with one stone. *Advances in neural information processing systems*, 35:22955–22968,
186 2022.

- 187 [17] Z. J. Cui, Y. Wang, N. Muhammad, L. Pinto, et al. From play to policy: Conditional behavior
188 generation from uncurated robot data. *arXiv preprint arXiv:2210.10047*, 2022.
- 189 [18] C. Chi, S. Feng, Y. Du, Z. Xu, E. Cousineau, B. Burchfiel, and S. Song. Diffusion policy:
190 Visuomotor policy learning via action diffusion. *arXiv preprint arXiv:2303.04137*, 2023.
- 191 [19] M. Reuss, M. Li, X. Jia, and R. Lioutikov. Goal-conditioned imitation learning using score-
192 based diffusion policies. *arXiv preprint arXiv:2304.02532*, 2023.
- 193 [20] S. Haldar, V. Mathur, D. Yarats, and L. Pinto. Watch and match: Supercharging imitation with
194 regularized optimal transport. *arXiv preprint arXiv:2206.15469*, 2022.
- 195 [21] S. Haldar, J. Pari, A. Rai, and L. Pinto. Teach a robot to fish: Versatile imitation from one
196 minute of demonstrations. *arXiv preprint arXiv:2303.01497*, 2023.
- 197 [22] A. Nair, A. Gupta, M. Dalal, and S. Levine. Awac: Accelerating online reinforcement learning
198 with offline datasets. *arXiv preprint arXiv:2006.09359*, 2020.
- 199 [23] B. Liu, Y. Zhu, C. Gao, Y. Feng, Q. Liu, Y. Zhu, and P. Stone. Libero: Benchmarking knowl-
200 edge transfer for lifelong robot learning. *arXiv preprint arXiv:2306.03310*, 2023.
- 201 [24] N. E. Sian, K. Yokoi, S. Kajita, F. Kanehiro, and K. Tanie. Whole body teleoperation of a
202 humanoid robot development of a simple master device using joysticks. *Journal of the Robotics
203 Society of Japan*, 22(4):519–527, 2004.
- 204 [25] Z. Gharaybeh, H. Chizeck, and A. Stewart. Telerobotic control in virtual reality. In *OCEANS
205 2019 MTS/IEEE SEATTLE*, pages 1–8, 2019. doi:10.23919/OCEANS40490.2019.8962616.
- 206 [26] T. Zhang, Z. McCarthy, O. Jow, D. Lee, X. Chen, K. Goldberg, and P. Abbeel. Deep imitation
207 learning for complex manipulation tasks from virtual reality teleoperation. In *ICRA*, 2018.
- 208 [27] S. P. Arunachalam, S. Silwal, B. Evans, and L. Pinto. Dexterous imitation made
209 easy: A learning-based framework for efficient dexterous manipulation. *arXiv preprint
210 arXiv:2203.13251*, 2022.
- 211 [28] I. Radosavovic, T. Xiao, S. James, P. Abbeel, J. Malik, and T. Darrell. Real-world robot learn-
212 ing with masked visual pre-training, 2022. URL <https://arxiv.org/abs/2210.03109>.
- 213 [29] A. George, A. Bartsch, and A. B. Farimani. Openvr: Teleoperation for manipulation. *arXiv
214 preprint arXiv:2305.09765*, 2023.
- 215 [30] A. G. Billard, S. Calinon, and F. Guenter. Discriminative and adaptive imitation in uni-manual
216 and bi-manual tasks. *Robotics and Autonomous Systems*, 54(5):370–384, 2006.
- 217 [31] A. Mandlekar, Y. Zhu, A. Garg, J. Booher, M. Spero, A. Tung, J. Gao, J. Emmons, A. Gupta,
218 E. Orbay, et al. Roboturk: A crowdsourcing platform for robotic skill learning through imita-
219 tion. In *Conference on Robot Learning*, pages 879–893. PMLR, 2018.
- 220 [32] P. Wu, Y. Shentu, Z. Yi, X. Lin, and P. Abbeel. Gello: A general, low-cost, and intuitive
221 teleoperation framework for robot manipulators. *arXiv preprint arXiv:2309.13037*, 2023.
- 222 [33] H. Fang, H.-S. Fang, Y. Wang, J. Ren, J. Chen, R. Zhang, W. Wang, and C. Lu. Low-cost ex-
223 oskeletons for learning whole-arm manipulation in the wild. *arXiv preprint arXiv:2309.14975*,
224 2023.
- 225 [34] A. Handa, K. Van Wyk, W. Yang, J. Liang, Y.-W. Chao, Q. Wan, S. Birchfield, N. Ratliff, and
226 D. Fox. Dexpivot: Vision-based teleoperation of dexterous robotic hand-arm system. In *2020
227 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9164–9170, 2020.
228 doi:10.1109/ICRA40945.2020.9197124.

- 229 [35] Y. Qin, W. Yang, B. Huang, K. Van Wyk, H. Su, X. Wang, Y.-W. Chao, and D. Fox.
230 Anytelep: A general vision-based dexterous robot arm-hand teleoperation system. *arXiv*
231 *preprint arXiv:2307.04577*, 2023.
- 232 [36] M. Caeiro-Rodríguez, I. Otero-González, F. A. Mikic-Fonte, and M. Llamas-Nistal. A system-
233 atic review of commercial smart gloves: Current status and applications. *Sensors*, 2021. ISSN
234 1424-8220. doi:10.3390/s21082667.
- 235 [37] V. Kumar and E. Todorov. Mujoco haptix: A virtual reality system for hand manipulation.
236 In *2015 IEEE-RAS 15th International Conference on Humanoid Robots (Humanoids)*, pages
237 657–663, 2015. doi:10.1109/HUMANOIDS.2015.7363441.
- 238 [38] S. Li, J. Jiang, P. Ruppel, H. Liang, X. Ma, N. Hendrich, F. Sun, and J. Zhang. A mobile robot
239 hand-arm teleoperation system by vision and imu. In *2020 IEEE/RSJ International Conference*
240 *on Intelligent Robots and Systems (IROS)*, pages 10900–10906. IEEE, 2020.
- 241 [39] S. Han, B. Liu, R. Cabezas, C. D. Twigg, P. Zhang, J. Petkau, T.-H. Yu, C.-J. Tai, M. Akbay,
242 Z. Wang, A. Nitzan, G. Dong, Y. Ye, L. Tao, C. Wan, and R. Wang. Megatrack: Monochrome
243 egocentric articulated hand-tracking for virtual reality. 2020.
- 244 [40] F. Zhang, V. Bazarevsky, A. Vakunov, A. Tkachenka, G. Sung, C.-L. Chang, and M. Grund-
245 mann. Mediapipe hands: On-device real-time hand tracking, 2020.
- 246 [41] D. A. Pomerleau. Alvin: An autonomous land vehicle in a neural network. In D. Touretzky,
247 editor, *NeurIPS*, volume 1. Morgan-Kaufmann, 1988.
- 248 [42] A. Y. Ng, S. J. Russell, et al. Algorithms for inverse reinforcement learning. In *ICML*, 2000.
- 249 [43] P. Abbeel and A. Y. Ng. Apprenticeship learning via inverse reinforcement learning. In *ICML*,
250 2004.
- 251 [44] I. Guzey, Y. Dai, B. Evans, S. Chintala, and L. Pinto. See to touch: Learning tactile dexterity
252 through visual incentives. *arXiv preprint arXiv:2309.12300*, 2023.
- 253 [45] D. A. Reynolds et al. Gaussian mixture models. *Encyclopedia of biometrics*, 741(659-663),
254 2009.
- 255 [46] S. Li, X. Ma, H. Liang, M. Görner, P. Ruppel, B. Fang, F. Sun, and J. Zhang. Vision-based
256 teleoperation of shadow dexterous hand using end-to-end deep neural network. In *2019 Inter-*
257 *national Conference on Robotics and Automation (ICRA)*, pages 416–422. IEEE, 2019.
- 258 [47] A. Sivakumar, K. Shaw, and D. Pathak. Robotic telekinesis: Learning a robotic hand imitator
259 by watching humans on youtube, 2022.
- 260 [48] Y. Qin, H. Su, and X. Wang. From one hand to multiple hands: Imitation learning for dexterous
261 manipulation from single-camera teleoperation. *arXiv preprint arXiv:2204.12490*, 2022.
- 262 [49] S. Li, N. Hendrich, H. Liang, P. Ruppel, C. Zhang, and J. Zhang. A dexterous hand-arm
263 teleoperation system based on hand pose estimation and active vision. *IEEE Transactions on*
264 *Cybernetics*, 2022.
- 265 [50] M. Mosbach, K. Moraw, and S. Behnke. Accelerating interactive human-like manipulation
266 learning with gpu-based simulation and high-quality demonstrations. In *2022 IEEE-RAS 21st*
267 *International Conference on Humanoid Robots (Humanoids)*, pages 435–441. IEEE, 2022.
- 268 [51] Z. Fu, T. Z. Zhao, and C. Finn. Mobile aloha: Learning bimanual mobile manipulation with
269 low-cost whole-body teleoperation. *arXiv preprint arXiv:2401.02117*, 2024.



271

Figure 2: We present OPEN TEACH, a unified robot teleoperation framework that supports multiple arms and hands, allows mobile manipulation, is calibration-free, and works across both simulation and real-world environments. OPEN TEACH uses a VR headset for teleoperation, offers low latency and high-frequency visual feedback. This high-frequency operation allows human users to correct for robot errors in real time, facilitating the execution of intricate and long-horizon tasks. From *making a sandwich* and *ironing cloth* to *placing items in a basket and lifting it* and *approaching a cabinet and opening it*, OPEN TEACH delivers a comprehensive, user-friendly teleoperation experience for a wide range of applications. OPEN TEACH is fully open-source.

272 **5.1 Framework details**

273 **5.1.1 Structure of the framework**

274 We use ZeroMQ for networking between nodes. The OPEN TEACH framework is divided into two
275 parts - *teleoperation* and *data collection*.

276 **Teleoperation:** The teleoperator is divided into 5 components - Detector, Keypoint Transformer,
277 Operator, Controller, and Visualizer. A brief description of each has been provided below.

- 278 1. **Detector:** Receives the hand keypoints from the Meta Quest 3 and publishes them to ZMQ
279 sockets.
- 280 2. **Keypoint Transformer:** Subscribes the keypoints published by the detector and maps
281 them to the robot pose.
- 282 3. **Operator:** Receives the robot pose from the keypoint transformer and the current robot
283 state from the controller. The operator computes the robot's actions which are published to
284 a ZMQ socket.

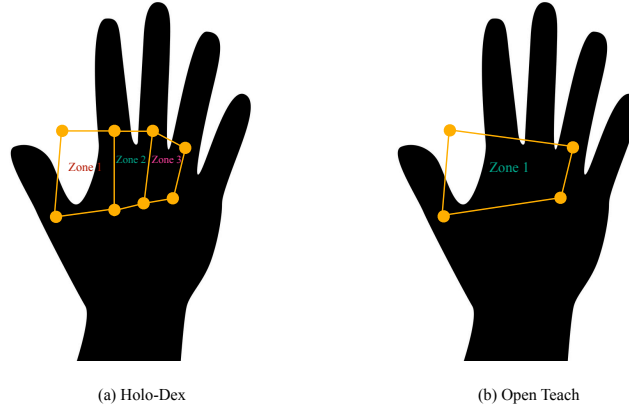


Figure 3: Thumb retargeting difference

- 285 4. **Controller:** Subscribes an action from the operator and takes an action in the real or simu-
 286 lated environment. After taking the action, the controller publishes the current states of the
 287 environment for use by the operator.
- 288 5. **Visualizer:** Subscribes the RGB images from the camera process (or the environment in
 289 case of simulations) and puts it on the screen inside the VR app for visualization during
 290 teleoperation.

291 **Data Collection:** A data recorder process subscribes sensor information (RGB and Depth images,
 292 tactile readings, timestamps) and robot-specific information (joint states, gripper states, timestamps)
 293 from the corresponding sockets and logs them in corresponding files. The data is then compiled
 294 together by matching the timestamps between the sensor information and robot-specific data.

295 5.1.2 Thumb Retargeting for Robot Hand

296 Section 2.3 provides details about the design of the OPEN TEACH wrapper for the robot hand. To
 297 recap, given the individual joint angles in the teacher’s hand from the VR headset, the joint angles
 298 for the robot hand can be computed by directly commanding the robot’s joints to the corresponding
 299 angles. This works well for all fingers except the thumb. Holo-Dex[13] deals with this by mapping
 300 the spatial coordinate of the teacher’s thumb tip to that of the robot hand. Then an inverse kinematics
 301 solver is used to compute the joint angles of the thumb. In this case, the retargeting of the thumb
 302 is done in 2D space. These bounds, depicted in Fig. 3(a), define the thumb’s reach limits. During
 303 retargeting, the thumb tip’s zone on the 2D palm plane is detected, and a perspective transform from
 304 the human hand to the robot hand is applied, aligning the human thumb tip with the robot thumb
 305 tip on the 2D plane. However, using three separate bounds introduces jitters when the thumb tip
 306 transitions between zones and results in stagnancy when outside the bounds. Further, in Holo-Dex,
 307 the height of the robot thumb tip is fixed, allowing it to only move along the 2D space.
 308 To address these challenges, OPEN TEACH employs a single, large zone spanning the entire thumb’s
 309 workspace in 2D space(refer to Fig. 3(b)). When the thumb is within bounds, a perspective trans-
 310 formation retargets the human thumb tip to the robot thumb tip. In cases where the thumb goes out
 311 of bounds, the closest point within the bound is estimated and used for retargeting, avoiding stagna-
 312 tion. Additionally, instead of a fixed height, the thumb is allowed to move perpendicular to the 2D
 313 surface along the palm, mapping the height of the human thumb tip to the robot thumb tip based on
 314 maximum and minimum height bounds. This approach ensures smoother thumb motion and enables
 315 the performance of more complex tasks compared to Holo-Dex [13].

316 6 Baseline Comparisons

317 Table 4 provides a comparison between OPEN TEACH and prior teleoperation systems considering
 318 features such as being calibration-free, compatible with multi-fingered hands, bimanual arms, and
 319 mobile manipulators, and being open-sourced.

Table 4: Comparison of OPEN TEACH’s capabilities with prior teleoperation systems on features such as being calibration-free, compatible with multi-fingered hands, bimanual arms, and mobile manipulators, and being open-sourced.

| | Calibration Free | Hands | Arms | Bimanual | Mobile Manipulation | Open-source |
|-------------------------|------------------|-------|------|----------|---------------------|-------------|
| Joystick | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ |
| Spacemouse | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ |
| Phone Teloperation [31] | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ |
| DexPilot [34] | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ |
| Holo-Dex [13] | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ |
| DIME [27] | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| TeachNet [46] | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ |
| Telekinesis [47] | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ |
| Qin et al. [48] | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ |
| MVP-Real [28] | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ |
| Transteleop [49] | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ |
| Mosbach et al. [50] | ✗ | ✓ | ✓ | ✗ | ✗ | ✓ |
| AnyTeleop [35] | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ |
| ALOHA [2] | ✓ | ✗ | ✓ | ✓ | ✗ | ✓ |
| Mobile ALOHA [51] | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ |
| GELLO [32] | ✓ | ✗ | ✓ | ✓ | ✗ | ✓ |
| AirExo [33] | ✓ | ✗ | ✓ | ✓ | ✗ | ✓ |
| Dobb-E [4] | ✓ | ✗ | ✓ | ✗ | ✓ | ✓ |
| OPEN TEACH | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

320 7 Experimental Setup

321 We evaluate the versatility of OPEN TEACH by using it to collect demonstrations on six different
 322 setups — four in the real world and two in simulation. Each setup is a combination of a variant of a
 323 robot arm with either an Allegro Hand or a 2-fingered gripper. The real-world setups include:

- 324 1. **Franka-Allegro:** A Franka Arm with an Allegro Hand having the Xela tactile sensors.
- 325 2. **Kinova-Allegro:** A Kinova Jaco Arm with an Allegro Hand with the Xela tactile sensors.
- 326 3. **Bimanual:** 2 xArm7 robot arms with 2-fingered grippers.
- 327 4. **Stretch:** Hello Stretch mobile manipulator with a 2-fingered gripper.

328 The Franka-Allegro and Kinova-Allegro comprise a single Intel Realsense camera for data collec-
 329 tion, whereas the Bimanual setup collects data from 5 different cameras. The Stretch has an iPhone
 330 attached to the wrist for data collection [4]. The simulated environments include:

- 331 1. **Allegro Sim:** A floating Allegro Hand capable of performing static and dynamic tasks.
- 332 2. **LIBERO Sim [23]:** A Franka Arm with a 2-fingered gripper placed in varied scenes.

333 8 Task Details

334 8.1 Demo Collections times

335 Table 5 provides the average times required to collect a demonstration for 16 tasks across 3 real-
 336 world setups (Franka-Allegro, Kinova-Allegro, Bimanual) and 2 simulated environments(Allegro
 337 sim, LIBERO sim).

338 8.2 Task Descriptions

339 Fig. 4, Fig. 5, Fig. 6, Fig. 7, Fig. 8, and Fig. 9 provide rollouts of all the tasks performed both in the
 340 real world and in simulated environments. Each task rollout is labeled with the name of the task and
 341 a task description.

342 8.3 Task Details

343 8.3.1 Demo Collections times

344 Table 5 provides the average times required to collect a demonstration for 16 tasks across 3 real-
 345 world setups (Franka-Allegro, Kinova-Allegro, Bimanual) and 2 simulated environments(Allegro
 346 sim, LIBERO sim).

Table 5: Time

| Robot Setup | Task | Average time to collect a demo (in s) |
|----------------|--|---------------------------------------|
| Franka-Allegro | Open box | 45 |
| | Grasp sponge | 60 |
| | Pick up tea satchet | 60 |
| | Grasp object and twist | 35 |
| Kinova-Allegro | Unfold towel | 40 |
| | Open a pack of cream | 10 |
| | Open ketchup bottle | 40 |
| Bimanual | Uncap marker | 60 |
| | Sweep table | 60 |
| | Pour sprinkles in a bowl | 40 |
| Allegro Sim | Flip cube | 3 |
| | Flip sponge | 20 |
| | Pinch Grasp | 15 |
| LIBERO Sim | Close top drawer of cabinet | 10 |
| | Turn on stove | 25 |
| | Pick up and put soup can in the basket | 30 |

347 8.3.2 Task Descriptions

348 Fig. 4, Fig. 5, Fig. 6, Fig. 7, Fig. 8, and Fig. 9 provide rollouts of all the tasks performed both in the
 349 real world and in simulated environments. Each task rollout is labeled with the name of the task and
 350 a task description.

351 8.4 User Study

352 Following up from Section ??, we provide the success rate and average completion times for all 15
 353 users for each task performed in Table 6 and Table 7 respectively. Each user roughly performed 3
 354 tasks on average, with 5 trials for each task. As mentioned in Section ??, since the Holo-Dex [13]
 355 and AnyTeleop [35] baselines lack open-source code for arm retargeting, we were unable to evaluate
 356 them on tasks involving arm movements. We observe a wide range of differences in success rates
 357 and average completion times demonstrating the inherent variations across users.

Table 6: Success rates for the user study conducted across 15 individuals. Each user roughly performs 3 tasks on average.

| User | Method | Success Rate (in 5 trials) | | | | |
|---------|------------|----------------------------|-------------|------|----------------|-------------------|
| | | Flip Cube | Pinch Grasp | Pour | Pick and Place | Open Box of Mints |
| User 1 | Holo-Dex | 1 | 0 | - | - | - |
| | AnyTeleop | 0.8 | 0.2 | - | - | - |
| | Open Teach | 1 | 0.8 | 0.2 | - | - |
| User 2 | Holo-Dex | - | 0.2 | - | - | - |
| | AnyTeleop | - | 0.2 | - | - | - |
| | Open Teach | - | 0.8 | - | 0.8 | 0.8 |
| User 3 | Holo-Dex | 1 | 0 | - | - | - |
| | AnyTeleop | 1 | 0.2 | - | - | - |
| | Open Teach | 1 | 0.8 | - | - | 0.2 |
| User 4 | Holo-Dex | 1 | 0 | - | - | - |
| | AnyTeleop | 1 | 0.2 | - | - | - |
| | Open Teach | 1 | 0.8 | - | 0.6 | 0.4 |
| User 5 | Holo-Dex | - | 0 | - | - | - |
| | AnyTeleop | - | 0.6 | - | - | - |
| | Open Teach | - | 0.2 | 0.4 | 1 | - |
| User 6 | Holo-Dex | - | 0 | - | - | - |
| | AnyTeleop | - | 0.6 | - | - | - |
| | Open Teach | - | 0.8 | - | 0.2 | - |
| User 7 | Holo-Dex | - | 0 | - | - | - |
| | AnyTeleop | - | 0 | - | - | - |
| | Open Teach | - | 0.6 | 0.8 | 0.8 | 0.4 |
| User 8 | Holo-Dex | 1 | - | - | - | - |
| | AnyTeleop | 1 | - | - | - | - |
| | Open Teach | 1 | - | - | - | - |
| User 9 | Holo-Dex | - | 0 | - | - | - |
| | AnyTeleop | - | 0.4 | - | - | - |
| | Open Teach | - | 0.8 | 0 | - | 0.6 |
| User 10 | Holo-Dex | - | 0 | - | - | - |
| | AnyTeleop | - | 0.2 | - | - | - |
| | Open Teach | - | 0.6 | 0.4 | 1 | 1 |
| User 11 | Holo-Dex | 1 | - | - | - | - |
| | AnyTeleop | 1 | - | - | - | - |
| | Open Teach | 1 | - | - | 0.8 | 0.4 |
| User 12 | Holo-Dex | 1 | - | - | - | - |
| | AnyTeleop | 1 | - | - | - | - |
| | Open Teach | 1 | - | - | - | - |
| User 13 | Holo-Dex | 1 | - | - | - | - |
| | AnyTeleop | 1 | - | - | - | - |
| | Open Teach | 1 | - | 0.6 | - | - |
| User 14 | Holo-Dex | - | 0 | - | - | - |
| | AnyTeleop | - | 0.4 | - | - | - |
| | Open Teach | - | 0.6 | - | - | 0.8 |
| User 15 | Holo-Dex | 1 | - | - | - | - |
| | AnyTeleop | 1 | - | - | - | - |
| | Open Teach | 1 | - | 0.4 | - | - |

Table 7: Average completion times for successful trials for the user study conducted across 15 individuals. Each user roughly performs 3 tasks on average. *NS* denotes cases where no successes were achieved.

| User | Method | Average completion time for successful demonstrations (in s) | | | | |
|---------|------------|--|-------------|-------|----------------|-------------------|
| | | Flip Cube | Pinch Grasp | Pour | Pick and Place | Open Box of Mints |
| User 1 | Holo-Dex | 4.6 | NS | - | - | - |
| | AnyTeleop | 20.2 | 22.5 | - | - | - |
| | Open Teach | 5.4 | 18.6 | 66 | - | - |
| User 2 | Holo-Dex | - | 17.5 | - | - | - |
| | AnyTeleop | - | 18.9 | - | - | - |
| | Open Teach | - | 20.6 | - | 29.7 | 12.2 |
| User 3 | Holo-Dex | 5.4 | NS | - | - | - |
| | AnyTeleop | 18.3 | 7.8 | - | - | - |
| | Open Teach | 5.1 | 12.6 | - | - | 11.3 |
| User 4 | Holo-Dex | 11 | NS | - | - | - |
| | AnyTeleop | 13.2 | 31.4 | - | - | - |
| | Open Teach | 6.2 | 7.5 | - | 16.9 | 48.4 |
| User 5 | Holo-Dex | - | NS | - | - | - |
| | AnyTeleop | - | 11.4 | - | - | - |
| | Open Teach | - | 10.9 | 41.6 | 12.4 | - |
| User 6 | Holo-Dex | - | NS | - | - | - |
| | AnyTeleop | - | 12.7 | - | - | - |
| | Open Teach | - | 10.5 | - | 23.57 | - |
| User 7 | Holo-Dex | - | NS | - | - | - |
| | AnyTeleop | - | NS | - | - | - |
| | Open Teach | - | 19.1 | 21.49 | 49 | 37.8 |
| User 8 | Holo-Dex | 6.5 | - | - | - | - |
| | AnyTeleop | 5.4 | - | - | - | - |
| | Open Teach | 4.7 | - | - | - | - |
| User 9 | Holo-Dex | - | NS | - | - | - |
| | AnyTeleop | - | 49.9 | - | - | - |
| | Open Teach | - | 65.3 | NS | - | 32.21 |
| User 10 | Holo-Dex | - | NS | - | - | - |
| | AnyTeleop | - | 48 | - | - | - |
| | Open Teach | - | 30.8 | 40.3 | 48.7 | 21.3 |
| User 11 | Holo-Dex | 6.7 | - | - | - | - |
| | AnyTeleop | 11.5 | - | - | - | - |
| | Open Teach | 5.6 | - | - | 21.8 | 15.7 |
| User 12 | Holo-Dex | 6.2 | - | - | - | - |
| | AnyTeleop | 11 | - | - | - | - |
| | Open Teach | 3.8 | - | - | - | - |
| User 13 | Holo-Dex | 8.9 | - | - | - | - |
| | AnyTeleop | 14.2 | - | - | - | - |
| | Open Teach | 5.8 | - | 18.1 | - | - |
| User 14 | Holo-Dex | - | NS | - | - | - |
| | AnyTeleop | - | 49.9 | - | - | - |
| | Open Teach | - | 65.3 | - | - | 132.5 |
| User 15 | Holo-Dex | 13.2 | - | - | - | - |
| | AnyTeleop | 14.6 | - | - | - | - |
| | Open Teach | 6.3 | - | 53.1 | - | - |

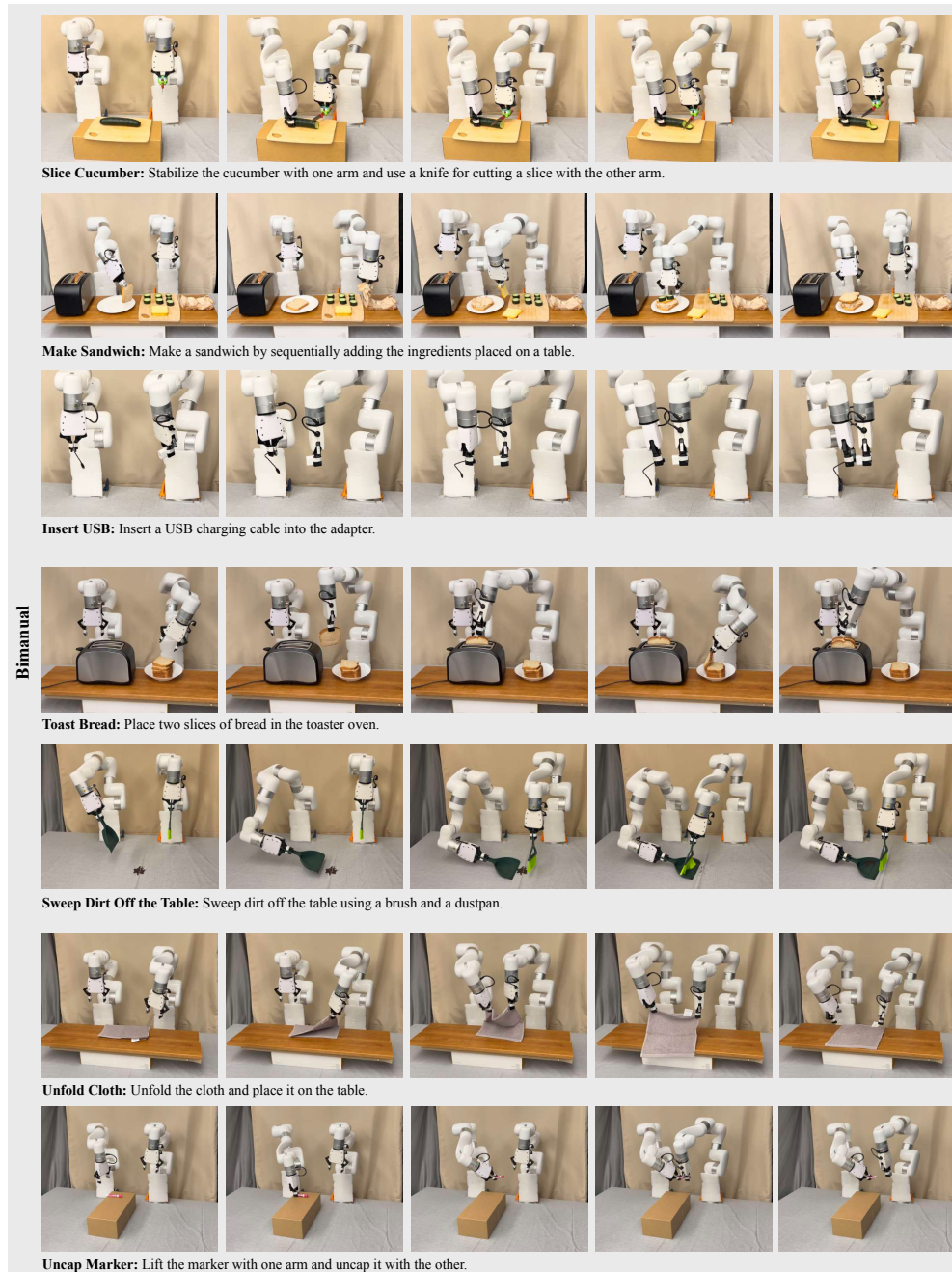


Figure 4: Real world task rollouts demonstrating the ability of OPEN TEACH to perform intricate, long-horizon tasks.

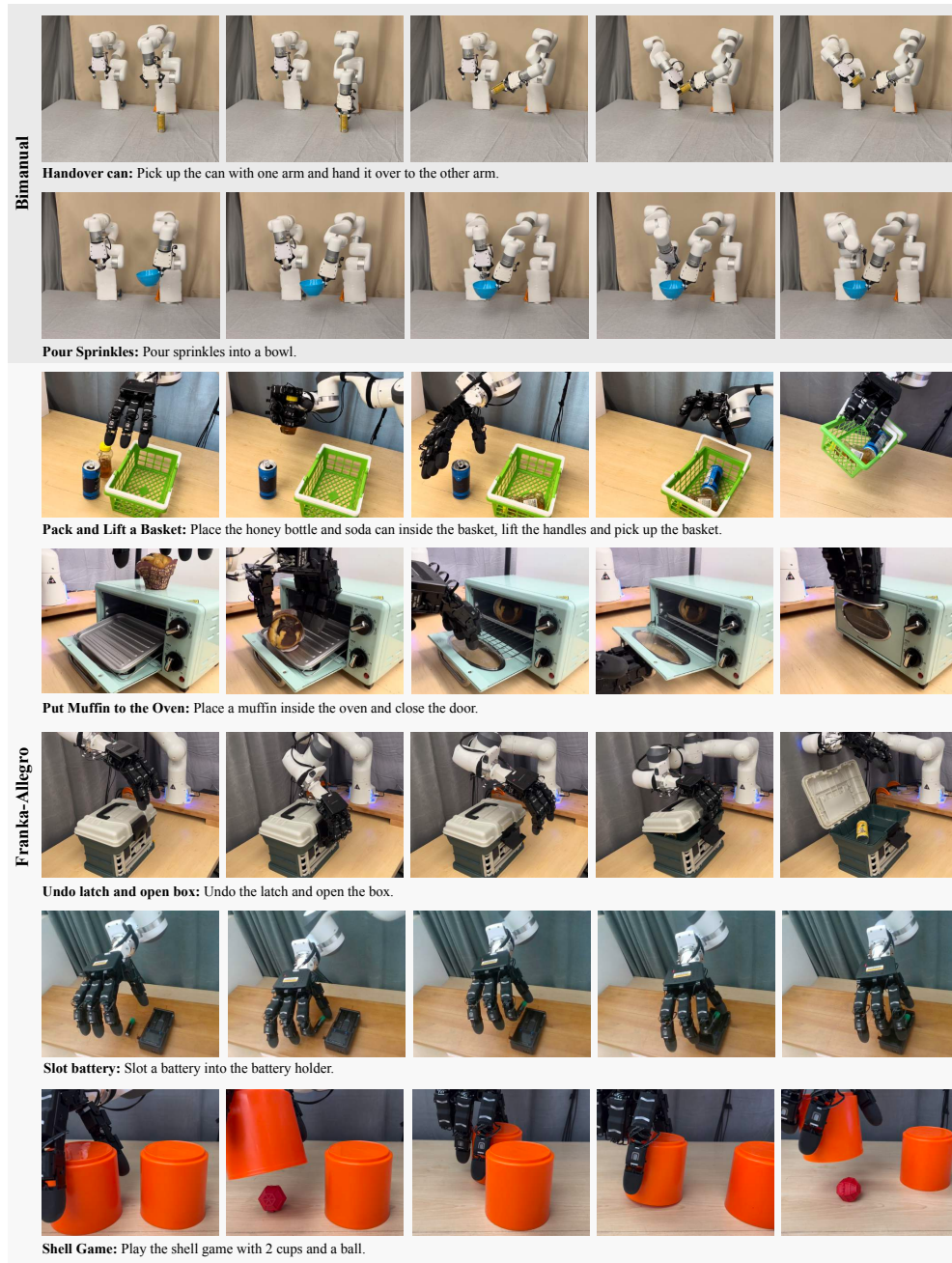


Figure 5: Real world task rollouts demonstrating the ability of OPEN TEACH to perform intricate, long-horizon tasks.



Figure 6: Real world task rollouts demonstrating the ability of OPEN TEACH to perform intricate, long-horizon tasks.

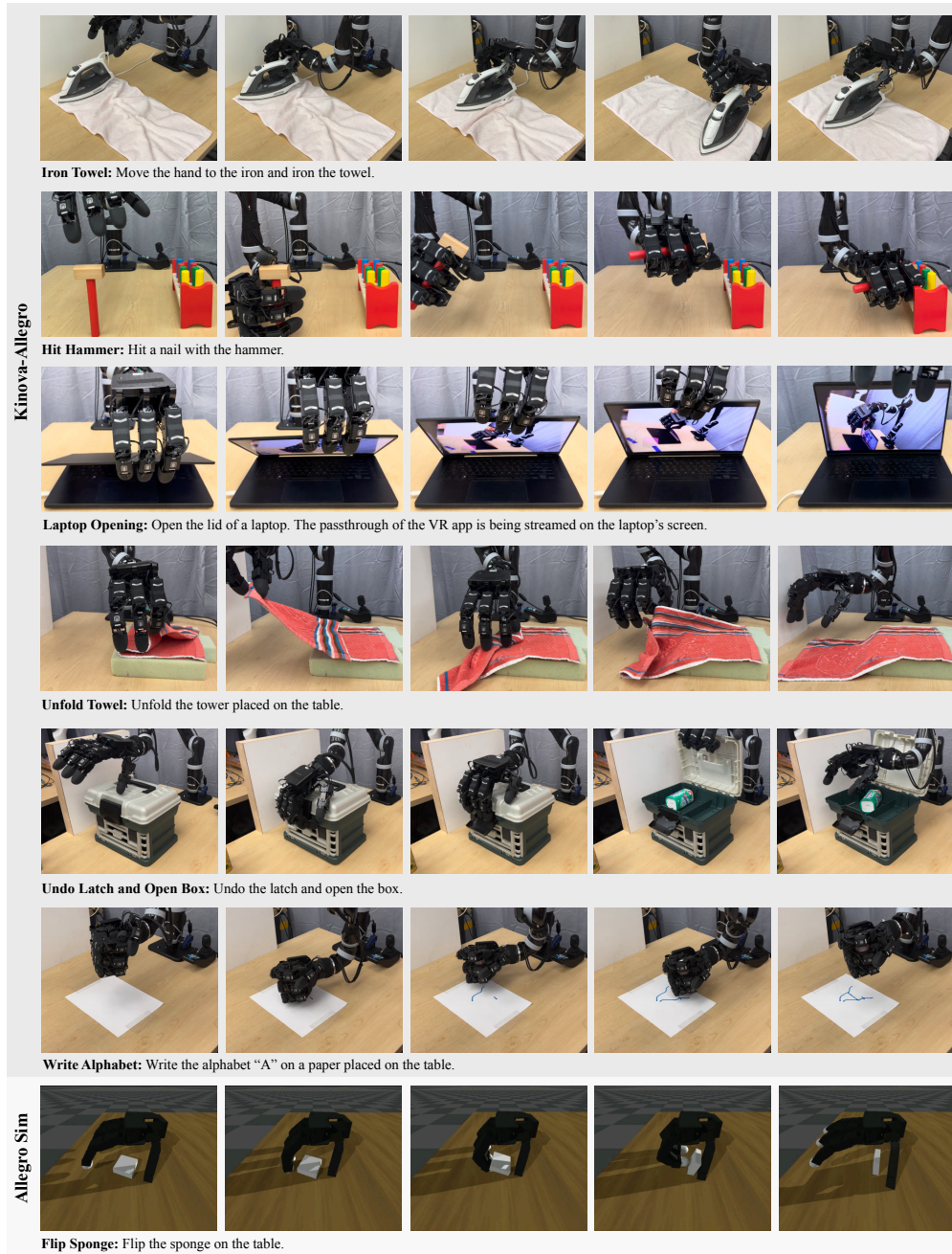


Figure 7: Real world task rollouts demonstrating the ability of OPEN TEACH to perform intricate, long-horizon tasks.

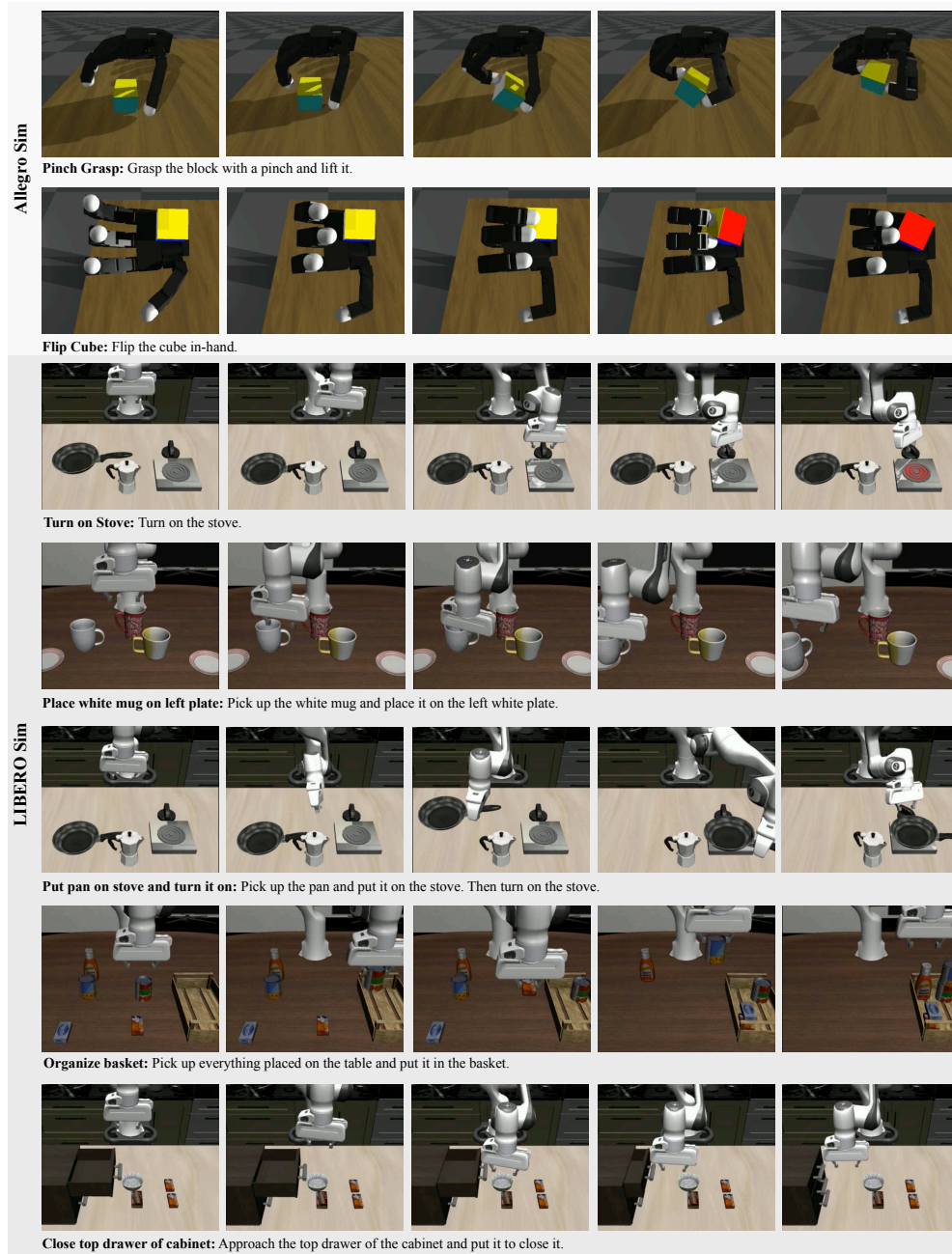


Figure 8: Real world task rollouts demonstrating the ability of OPEN TEACH to perform intricate, long-horizon tasks.



Figure 9: Real world task rollouts demonstrating the ability of OPEN TEACH to perform intricate, long-horizon tasks.

Table 8: Performance of policies learned on data collected through OPEN TEACH. FISH and BC were used to train policies for Allegro Sim and Libero Sim respectively. We report the mean and standard deviation for 25 evaluation trials across 3 seeds for each task.

| Robot Setup | Task | Number of Demos | Success Rate (25 trials) |
|-------------|---------------------------------|-----------------|--------------------------|
| Allegro Sim | Flip Cube | 6 | 0.97 ± 0.03 |
| | Flip Sponge | 6 | 0.79 ± 0.05 |
| | Pinch Grasp | 6 | 0.75 ± 0.07 |
| Libero Sim | Close Top Drawer of Cabinet | 10 | 0.96 ± 0.03 |
| | Turn on Stove | 10 | 0.95 ± 0.04 |
| | Pick and Place Soup into Basket | 50 | 0.77 ± 0.02 |

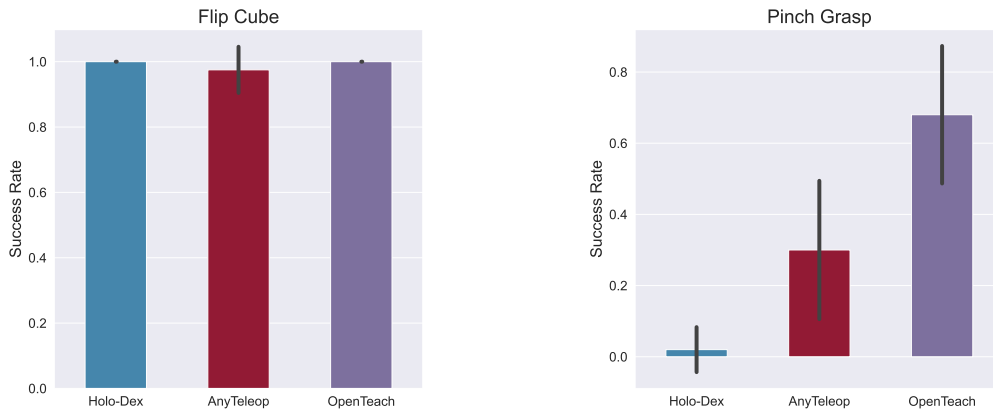


Figure 10: Success rates for the user study conducted across 15 individuals on 2 tasks - Flip Cube and Pinch Grasp. We report the mean and standard deviation for 3 methods - Holo-Dex, AnyTeleop, and Open Teach (Ours).

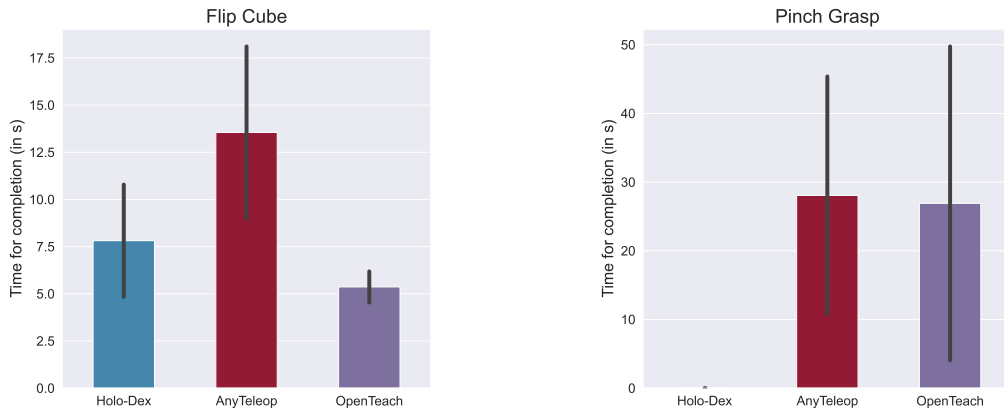


Figure 11: Average completion times for successful trials for the user study conducted across 15 individuals for 2 tasks - Flip Cube and Pinch Grasp. We report the mean and standard deviation for 3 methods - Holo-Dex, AnyTeleop, and Open Teach (Ours).

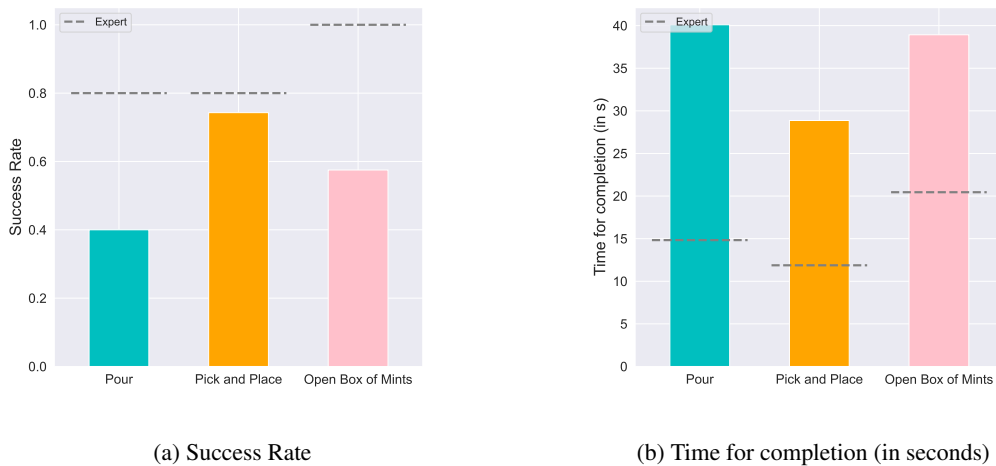


Figure 12: We compare the (a) success rate, and (b) average completion time (in seconds) for using OPEN TEACH between an expert and 15 individuals participating in a user study. We report this comparison for 3 tasks - Pour, Pick and Place, and Open Box of Mints.