# OPEN TEACH: A Versatile Teleoperation System for Robotic Manipulation

Aadhithya Iyer[1]    Zhuoran Peng[1]    Yinlong Dai[1]    Irmak Güzey[1]
Siddhant Haldar[1]    Soumith Chintala[2]    Lerrel Pinto[1]

[1]New York University    [2]Meta

https://open-teach.github.io/

Fig. 1: 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 and 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, 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.

*Abstract*—**Open-sourced, user-friendly tools form the bedrock of scientific advancement across disciplines. The widespread adoption of data-driven learning has led to remarkable progress in multi-fingered dexterity, bimanual manipulation, and applications ranging from logistics to home robotics. However, existing data collection platforms are often proprietary, costly, or tailored to specific robotic morphologies. We present OPEN TEACH, a new teleoperation system leveraging VR headsets to immerse users in mixed reality for intuitive robot control. Built on the affordable Meta Quest 3, which costs $500, OPEN TEACH enables real-time control of various robots, including multi-fingered hands, bimanual arms, and mobile manipulators, through an easy-to-use app. Using natural hand gestures and movements, users can manipulate robots at up to 90Hz with smooth visual feedback and interface widgets offering closeup environment views. We demonstrate the versatility of OPEN TEACH across 38 tasks on different robots. A comprehensive user study indicates significant improvement in teleoperation capability over the AnyTeleop framework. Further experiments exhibit that the collected data is compatible with policy learning on 10 dexterous and contact-rich manipulation tasks. Currently supporting Franka, xArm, Jaco, Allegro, and Hello Stretch platforms, OPEN TEACH is fully open-sourced to promote broader adoption. Videos are available at https://open-teach.github.io/.**

## I. INTRODUCTION

The integration of learning-based methods has sparked a revolution in robotics, significantly enhancing capabilities in manipulation [10, 67, 7, 53], locomotion [17, 36, 57, 9], and aerial robotics [65, 16, 26]. More recent work has been making advancements in complex single-task behavior learning [60, 4, 67], multitask scenarios [43, 5], multimodal behavior learning [52, 13, 11, 49], and efficient fine-tuning of pretrained behavior models [21, 22, 40]. A fundamental requirement across all these threads of research is the need to collect data in the form of task demonstrations.

Commonly used teleoperation systems include devices such as joysticks and 3D spacemouses [34, 55], commercial VR headsets [19, 66, 3, 4, 48, 18], kinesthetic teaching [6], and phone teleoperation [37]. The aforementioned devices are cost-effective and easy to set up. However, they are often unintuitive to use and require extensive user-training to demonstrate intricate motions. Recently proposed exoskeleton-based teleoperation frameworks like ALOHA [67], GELLO [61], and AirExo [14] attempt to alleviate this problem by having the human teleoperator directly control a kinematically isomorphic version of the same robot arm. These frameworks directly impose the kinematic constraints of the robot arm during teleoperation making it more compatible and intuitive to control the motion of the robot. Although highly effective, these systems can require an additional robot for each robot being controlled, have high initial setup costs, and are designed for specific robot morphologies.

The challenge of easy-to-use teleoperation devices is more apparent in dexterous manipulation problems [24, 47, 3, 4], owing to the high dimensional action space. Such frameworks typically involve the use of expensive gloves [8, 29, 32],

Correspondence to: aadhithya.iyer@nyu.edu

extensive calibration processes [24, 3], or are susceptible to monocular occlusions [3].

In this work, we present OPEN TEACH, an open-source framework for robot teleoperation that supports a variety of robots, including bimanual and multi-finger manipulation, all at a price of $500. As shown in Figure 1, OPEN TEACH uses a VR headset (e.g. Quest 3) to put users / teachers in an immersive virtual world where they can view a robotic scene both through their eyes, via visual passthrough, as well as realtime streams from the robot's cameras. To control the robot, users can simply use hand gestures, which are detected using onboard hand-pose estimators at 90Hz. As a result, even though OPEN TEACH is kinematics-unaware, the high frequency execution and improved hand pose detection accuracy enables users to collect high-quality real-time robot demonstrations.

We experimentally evaluate OPEN TEACH on 38 tasks across single arm, bimanual, multi-fingered, and mobile manipulation robot setups in both simulation and the real world. The tasks span from tabletop manipulation to contact-rich dexterous manipulation. Across different robot morphologies, we find that users can provide demonstrations at speeds on par with robot-specific teleoperation systems and significantly faster than general-purpose systems like AnyTeleop [47]. Importantly, policies trained on the data collected achieve an average success rate of 86% on 10 tasks in simulation and the real world, validating the utility of policy learning using OPEN TEACH. The contributions of this work is summarized as follows:

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

OPEN TEACH is fully open-source with the mixed reality API, policy training code, and demonstrations collected using OPEN TEACH available at https://open-teach.github.io/.

## II. RELATED WORK

### A. Robot-Specific Teleoperation

Teleoperation, as a medium for human-robot interaction, has been a crucial part of robotics. Recent strides in learning-based methods demand extensive data collection, giving rise to diverse teleoperation systems — joysticks and spacemouses [34, 55], VR controllers [19, 66, 3, 4, 48, 18], RGB cameras [24, 56, 47, 58], IMU sensors [30, 28, 62], kinesthetic teaching [6], phone teleoperation [37], gloves [8, 29, 32], marker-based motion capture systems [68, 35], and reacher-grabber sticks [44, 53]. However, several challenges remain in the robot-specific nature of these devices. For instance, devices like joysticks, spacemouses, and phones are limited to

TABLE I: 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 Teleoperation [37] | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ |
| DexPilot [24] | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ |
| Holo-Dex [4] | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ |
| DIME [3] | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| TeachNet [31] | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ |
| Telekinesis [56] | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ |
| Qin et al. [46] | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ |
| MVP-Real [48] | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ |
| Transteleop [33] | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ |
| Mosbach et al. [39] | ✗ | ✓ | ✓ | ✗ | ✗ | ✓ |
| AnyTeleop [47] | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ |
| ALOHA [67] | ✓ | ✗ | ✓ | ✓ | ✗ | ✓ |
| Mobile ALOHA [15] | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ |
| GELLO [61] | ✓ | ✗ | ✓ | ✓ | ✗ | ✓ |
| AirExo [14] | ✓ | ✗ | ✓ | ✓ | ✗ | ✓ |
| Dobb-E [53] | ✓ | ✗ | ✓ | ✗ | ✓ | ✓ |
| OPEN TEACH | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

controlling robot arms due to their lack of fidelity for multi-fingered hands. There have been systems developed to map the human hand pose to the robot pose [3, 4, 47, 31, 2, 38], but they are often restricted to only controlling robot hands. Further, all of these frameworks lack awareness of the robot's kinematic constraints, leading to challenges in intuitive control, especially in complex poses. As a solution, there are more conventional but expensive exo-skeleton based teleoperation systems [25, 27, 51] that use a second arm for controlling the manipulator arm. Recently proposed ALOHA [67, 15] affirms the effectiveness of this approach through impressive results in fine-grained bimanual manipulation. However, these systems require an exact copy of the manipulator robot arm, rendering them costly and less practical for heavier robots. In addressing these issues, GELLO [61] and AirExo [14] introduce exo-skeleton teleoperation frameworks, utilizing a kinematically isomorphic variant of the robot arm. This approach proves more affordable and lightweight, enhancing usability for humans. Despite their success in fine-grained manipulation tasks, these solutions are constrained to robot arms and face challenges in extending seamlessly to control robot hands. For multi-fingered robot hands, gloves, vision-based, and VR-based methods have been employed. These systems either assume a fixed robot arm [3, 4] or are tied to a specific robot setup [24, 33], making them difficult to transfer to new arm-hand systems and new environments.

### B. Unified Teleoperation Frameworks

The robotics community has often sought to develop versatile systems that operate across diverse environments and robots [54, 43, 42]. Leveraging the success of learning-based methods, achieving this requires teleoperation systems adaptable to various robot variants, allowing for abundant data collection with minimal setup costs. While methods combining robot arms with multi-fingered hands exist [24, 56, 48, 39,

33], their applicability across robot variants remains unclear. AnyTeleop [47] makes progress in this direction by proposing a robot-agnostic system compatible with multiple hands and arms. We build upon this idea in creating OPEN TEACH.

## III. BACKGROUND ON IMITATION LEARNING

### A. Behavior Cloning

Given a dataset of expert rollouts for a desired task in the form of observation and action pairs $\mathcal{D} \equiv \{(o, a)\} \subset \mathcal{O} \times \mathcal{A}$, behavior cloning (BC) aims to learn a policy $\pi: \mathcal{O} \rightarrow \mathcal{A}$ that models this data without any online interactions with the environment nor a reward function. Often, such policies are chosen from a hypothesis class parameterized by a parameter set $\theta$. Following this convention, the objective of BC is to find the value $\theta$ that maximizes the probability of the observed data.

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \prod_t \mathbb{P}(a_t|o_t; \theta) \qquad (1)$$

When constrained to unimodal isotropic Gaussians, this maximum likelihood estimation problem leads to minimizing the Mean Squared Error (MSE), $\Sigma_t \|a_t - \pi(o_t; \theta)\|^2$.

### B. Inverse Reinforcement Learning

In this work, we employ FISH [22] and TAVI [20] to learn visual and visuotactile policies respectively. For both of these methods, the first phase involves obtaining a non-parametric base-policy $\pi^b: \mathcal{Z} \rightarrow \mathcal{A}$ with encoded representations $z \in \mathcal{Z}$ and actions $a \in \mathcal{A}$. Then a residual policy $\pi^r: \mathcal{Z} \times \mathcal{A} \rightarrow \mathcal{A}$ is learned atop the base policy $\pi^b$ such that an action sampled from the final policy $\pi$ is the sum of the base action $a^b \sim \pi^b(z)$ and the residual offset $a^r \sim \pi^r(z, a^b)$. The reward for learning the residual policy through inverse RL is obtained through optimal transport based trajectory matching [21, 12].
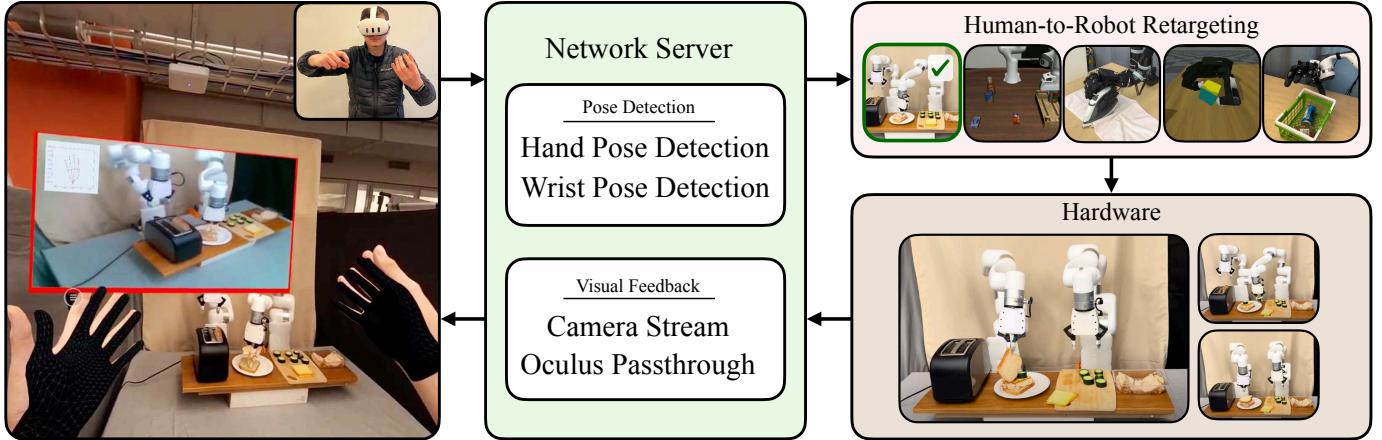
Fig. 2: Overview of the teleoperation module in OPEN TEACH. Provided a hand and wrist pose within the VR interface, the controller transmits keypoint data to the robot's server. The server then transforms and retargets these key points to align with the specific robot setup. Real-time visual feedback of the teleoperated scene is promptly relayed back to the VR headset.

## IV. OPEN TEACH

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

In this section, we provide details about the VR-based teleoperation setup and the system design that enables data collection using this framework.

### A. Placing an Operator in a Virtual World

We use the Meta Quest 3 VR headset to place the human teacher in a virtual world. The headset surrounds the human in a virtual environment at a resolution of $2064 \times 2208$ and a native refresh rate of 90Hz. The base version of this headset is affordable at $499 and is relatively light at 513g. Compared to the Meta Quest 2 VR headset used in prior work [4], the Quest 3 provides a full-color passthrough allowing the human to get a direct view of the robot setup during teleoperation. These features, especially the full-color passthrough, allow for a comfortable and intuitive operation by the user. Additionally, similar to Arunachalam et al. [4], the Quest 3 API interface allows for creating custom mixed reality worlds that visualize the robotic system along with diagnostic panels in VR. Examples of virtual scenes have been shown in Fig. 2 and Fig. 3. It is important to highlight the exceptional clarity of the scene passthrough visible in Quest 3. The teacher can experience a 3D view of the scene through the headset, significantly enhancing the intuitiveness of the teleoperation experience.

### B. Pose Estimation with VR Headsets

Similar to Arunachalam et al. [4], we directly use the in-built hand pose estimator [23] of the Quest 3 using 2 monochrome cameras. This is significantly more robust compared to single camera alternatives [64]. Further, since the cameras are internally calibrated, they do not require specialized calibration routines that are needed in prior multi-camera teleoperation frameworks [24, 47]. Also, since the hand-pose estimator is integrated into the device, it can stream real-time hand poses at 90Hz. This alleviates the challenge of obtaining hand poses at both high accuracy and high frequency, as reported in prior work [24, 3].

### C. Human to Robot Pose Retargeting

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

**Robot Arm**: We use the wrist keypoint and the knuckle points of the index and pinky fingers to establish a 3D coordinate system with a 2D plane along the palm of the human hand, and a third axis perpendicular to the palm. Then, the wrist position is mapped to the robot end effector position, and the transformations of the 3D coordinate system across time are mapped to the changes in the end effector orientation.

**Robot Hand**: We use the teacher's hand pose obtained from the VR to compute the individual joint angles in the teacher's hand. Given these joint angles, a straightforward method of retargeting is to directly command the robot's joints to the corresponding angles. In practice, this works well for all fingers except the thumb. To address this, we improve upon Arunachalam et al. [4], where the spatial coordinate of the teacher's thumb tip is mapped to that of the robot hand and then an inverse kinematics solver is used to compute the joint angles of the thumb. More details about the improvement in
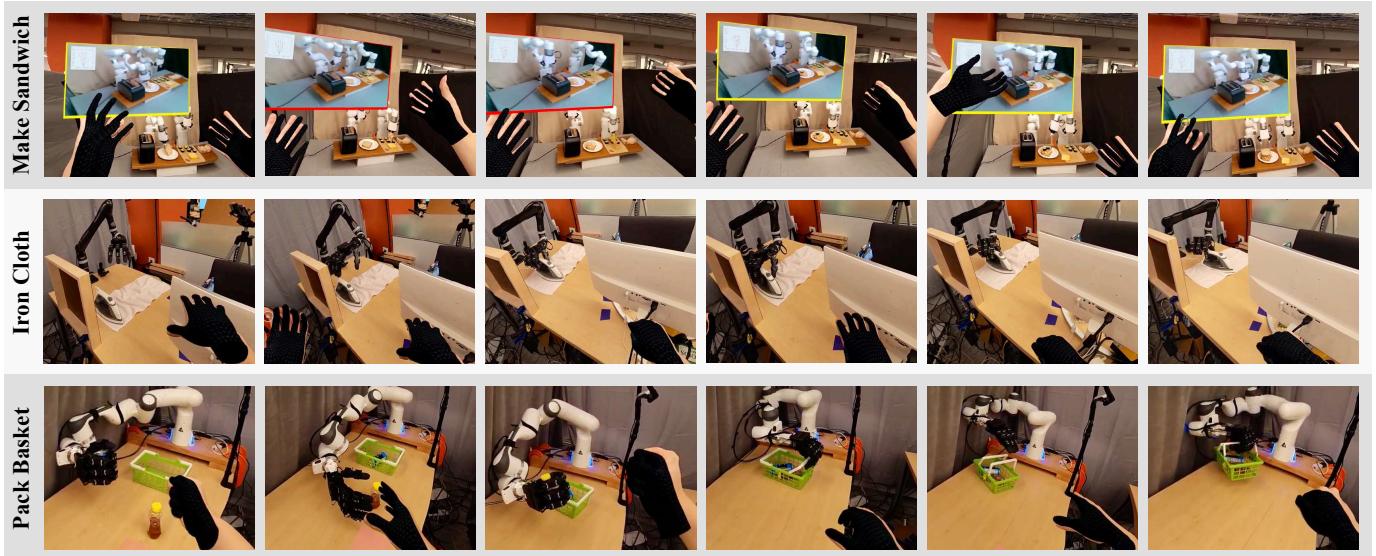
Fig. 3: The demonstration collection process as viewed from within the VR application. Shown here is one task being performed for each real-world setup. High resolution images are streamed at 90 Hz to the VR application, allowing for an immersive experience and improved teleoperation compared to prior VR-based systems. The high frame rate streaming enables reactive control by the user, while widgets for visualizing the robot's camera view help the user focus on fine-grained movements.

thumb retargeting have been included in Appendix A. Since the Allegro hand does not have a pinky finger, we ignore the teacher's pinky joints.

**Two-fingered gripper**: To detect the opening and closing of the two-fingered gripper, we utilize the pinch between the pinky finger and the thumb. The pinch is detected by computing the distance between the tips of the two fingers and setting a threshold on the pinch distance. We use a toggle mechanism for opening and closing the gripper where each pinch indicates toggling to the alternate state of the gripper.

**Mobile manipulator**: The same 3D coordinate system established for controlling robot arms is used for mapping the wrist's movements to actions of the mobile robot. When the wrist moves forward, it extends the robot's arm, enabling it to reach farther. Vertical wrist movements adjust the robot's height, while lateral wrist movements cause the robot to move sideways by controlling its wheels. The 3D transformations across time are mapped to the changes in the end effector orientation. The opening and closing of the gripper are controlled through the pinch between the index finger and the thumb.

These are just a few examples of how the hand pose data can be used to obtain a mapping between the human hand and a robot. The simplicity of the proposed framework has been summarized in Code Snippet 1. The primary idea behind OPEN TEACH is that given a VR headset, the end-user has the flexibility to design their own human-to-robot wrappers using the human hand poses as input. The framework has been designed for simple integration with any robot setup, allowing robot teleoperation with real-time streaming (up to 90Hz) and low-latency visual feedback. This significantly reduces the initial setup cost as compared to prior exoskeleton-based teleoperation frameworks like GELLO [61] and AirExo [14].

**Code Snippet 1:** Robot control using VR

```
# Initialize robot and VR
vr = VR()
robot = Robot()
while(True):
    # Step 1: Get hand pose from VR
    hPose = vr.getHandPose()
    # Step 2: Retarget to robot pose
    rPose = robot.retargetH2R(hPose)
    # Step 3: Move robot to target pose
    robot.move(rPose)
```

To increase support for more robots, we invite roboticists to add support for their robots on our common OPEN TEACH GitHub repository (see Section VII).

#### D. Robot Control

Achieving minimal error and low latency is pivotal for OPEN TEACH to facilitate the intuitive teleoperation of the robot hand by the human teacher. In this study, we employ the Allegro Hand as our robotic hand, controlling it asynchronously through the ROS [59] communication framework. Using the computed robot joint positions from the retargeting procedure, a PD controller outputs desired torques at a frequency of 300Hz. To mitigate steady-state error, we include a gravity compensation module to compute offset torques.

We use three different robot arms for our evaluations — xArm, Franka Emika Panda, and Kinova Jaco. We use different controllers for each. The xArm is directly controlled through the official xArm API [63]. For the Franka Emika Panda, we use the Deoxys controller [69]. For the Kinova Jaco, we use

the controller open-sourced by Arunachalam et al. [3]. The Allegro hand's streaming frequency is configured at 60 Hz, while the xArm, Franka Emika Panka, and Kinova Jaco arms are set to 90 Hz, 60 Hz, and 60 Hz, respectively. Such high frequency teleoperation allows the human teacher to see the robot move in real time and immediately correct execution errors in the robot. The Hello Stretch is controlled at 5Hz using the controller released by Shafiullah et al. [53]. Further, we acknowledge the fact that the human hand possesses fewer degrees of freedom than the robot. In response, we introduce a pause functionality, allowing the teacher to momentarily halt teleoperation, reorient themselves, and resume the process. Furthermore, we implement a resolution adjustment feature, which provides a performance boost for high-precision tasks such as delicately picking up a tea sachet. Details about these implementations have been included in Appendix A.

## V. Experimental Evaluation

Our experiments and tasks are designed to answer the following questions:

1) How versatile is OPEN TEACH across a range of robotics setups?
2) How successful are policies trained with OPEN TEACH?
3) Can OPEN TEACH be used for performing complex, long-horizon tasks?
4) How intuitive is the system for new users?

### A. Experimental Setup

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

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

The Franka-Allegro and Kinova-Allegro comprise a single Intel Realsense camera for data collection, whereas the Bimanual setup collects the data from 5 different cameras. The Stretch has an iPhone attached to the wrist for data collection, similar to Shafiullah et al. [53]. The simulated environments include:

1) **Allegro Sim:** Comprising a floating Allegro Hand capable of performing both static and dynamic tasks.
2) **LIBERO Sim [34]:** Comprising a Franka Arm with a 2-fingered gripper placed in varied scenes.

We demonstrate the usefulness of the collected data by training visual and visuotactile policies using behavior cloning [45] and inverse RL [41, 1].

### B. Imitation Learning with OPEN TEACH Data

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

1) **Franka-Allegro:** We record both visual and tactile data for this setup. The policies are trained using TAVI [20], a demonstration-guided residual RL algorithm that collects a few expert demonstrations and learns a robot policy using both visual and tactile data.
2) **Allegro Sim:** We only record visual data for this setup. The policies are trained using FISH [22].
3) **LIBERO Sim [34]:** We only record visual data for this setup. The policies are trained using transformer-based BC with a GMM head [50] and action chunking [67].

### C. How versatile is OPEN TEACH across robotic setups?

The primary idea behind OPEN TEACH is that given any robotic setup, a user can purchase an affordable off-the-shelf VR headset (in this case, Quest 3) and plug the headset and robot setup into the proposed framework to start teleoperating the robot without any additional hardware setup cost. To investigate its versatility, we use OPEN TEACH to teleoperate four different real world robotic setups, each having a different combination of a robot arm and end effector type — Franka Allegro, Kinova Allegro, a Bimanual setup with 2 xArm7 robots, and Hello Stretch for mobile manipulation. We also exhibit the applicability of OPEN TEACH in simulation through evaluations on 2 simulated environment suites — Allegro Sim and LIBERO Sim [34]. The frequency of teleoperation for each of the setups has been provided in Table II. Table III provides a set of tasks performed on Franka-Allegro, Allegro Sim, and LIBERO Sim. A more comprehensive list of tasks, including those on the Kinova-Allegro, Bimanual, and Stretch setup have been provided in Fig. 4 and Appendix B2.

### D. How successful are policies trained with OPEN TEACH?

Table III provides the success rates of policies learned using imitation learning across both the real-world and simulated setups. We use TAVI [20] to learn visuotactile policies on Franka-Allegro, and FISH [22] to learn visual policies on Allegro Sim. Similar to prior work [20, 22], these policies were learned within 20 minutes and achieved an average success rate of 82%, validating the high quality of the observation data collected through OPEN TEACH. We employ behavior cloning to train policies on LIBERO Sim, achieving an average success rate of 93%, thus confirming the high quality of the collected action data. Overall, the learned policies achieve an average success rate of 86% across all tasks and robot morphologies. This underscores the effectiveness of OPEN TEACH in collecting data for policy learning.

### E. Can OPEN TEACH be used for performing complex, long-horizon tasks?

In this section, we emphasize the efficacy of OPEN TEACH in executing a diverse array of complex, long-horizon tasks across various robotic configurations. Illustrated in Fig. 4 are

**Make Sandwich:** Make a sandwich by sequentially adding the ingredients placed on a table.

**Insert USB:** Insert a USB charging cable into the adapter.

**Pack and Lift a Basket:** Place the honey bottle and soda can inside the basket, lift the handles and pick up the basket.

**Put Muffin to the Oven:** Place a muffin inside the oven and close the door.

**Iron Towel:** Move the hand to the iron and iron the towel.

**Laptop Opening:** Open the lid of a laptop. The passthrough of the VR app is being streamed on the laptop's screen.

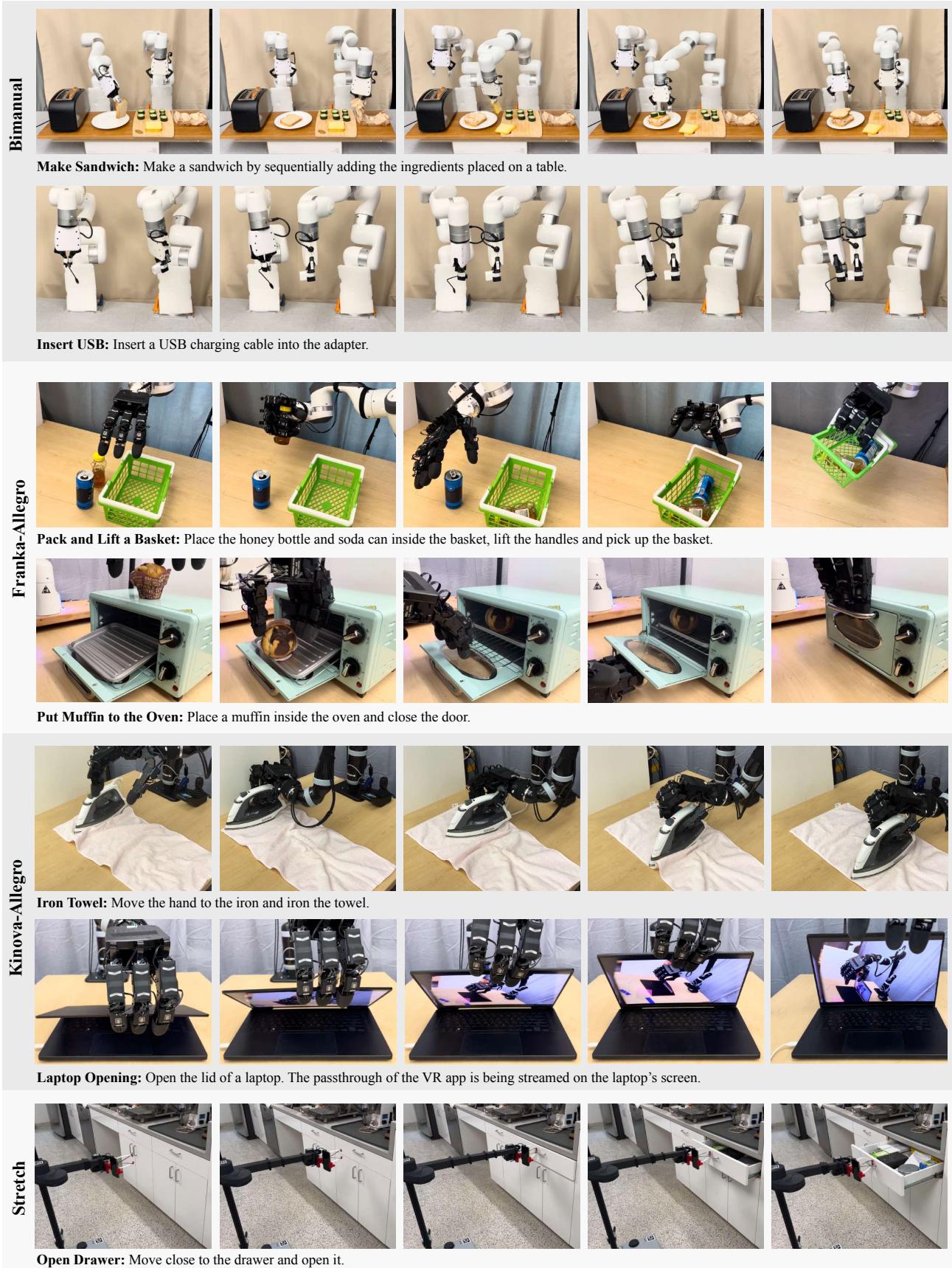**Open Drawer:** Move close to the drawer and open it.

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

TABLE II: Teleoperation Frequency across all robots.

| 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 |

TABLE III: Performance of policies learned on data collected through OPEN TEACH. For Franka-Allegro, Allegro Sim, and Libero Sim, TAVI [20], FISH [22] and BC were respectively used to train policies.

| 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 |
| | Grasp Object and Twist | 6 | 8/10 |
| Allegro Sim | Flip Cube | 6 | 10/10 |
| | Flip Sponge | 6 | 10/10 |
| | Pinch Grasp | 6 | 7/10 |
| Libero Sim | Close Top Drawer of Cabinet | 10 | 10/10 |
| | Turn on Stove | 10 | 9/10 |
| | Pick and Place Soup into Basket | 50 | 9/10 |

TABLE IV: User study comparing OPEN TEACH with other 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 |

examples of real-world task rollouts for the Bimanual, Franka-Allegro, Kinova-Allegro and Stretch setups. OPEN TEACH allows the collection of demonstrations for intricate, extended tasks, ranging from high-precision activities like USB insertion to delicate movements such as slicing a cucumber. On the multi-fingered hand setup, we demonstrate a broad spectrum of tasks, encompassing extended activities like *placing objects in a basket and lifting it* to contact-rich manipulation scenarios like *opening a laptop* and *sliding a tea sachet off the table*. A detailed compilation of tasks performed in both real-world and simulated setups, along with more task rollouts, have been included in Appendix B. Videos showcasing these task rollouts can be found on our project website.

*F. How intuitive is the system for new users?*

We assess the user-friendliness of the OPEN TEACH through a comprehensive user study involving 15 new users. The study is carried out on the Franka-Allegro setup, chosen for its capacity to evaluate the system's performance on both the robot hand and the robot arm. Each participant is allocated a 10-minute practice session to familiarize themselves with teleoperating the robot setup. Following this, they are tasked with performing five trials for each of three distinct tasks using Holo-Dex [4], AnyTeleop [47], and OPEN TEACH. To mitigate potential biases, the order of tasks is randomized for each user.

In Table IV, we present a comparative analysis of success rates and median completion times for new users across Holo-Dex, AnyTeleop, and OPEN TEACH for the tasks of cube flipping and pinch grasping. We chose to analyze the median rather than mean completion times in order to mitigate potential biases from outliers, given the relatively small user sample size. Since the Holo-Dex and AnyTeleop baselines

lack open-source code for arm retargeting, we were unable to evaluate them on tasks involving arm movements. Thus, our comparison is limited to the cube flipping and pinch grasping tasks that do not require arm manipulation. On these tasks, OPEN TEACH demonstrates a higher success rate along with significantly reduced median time to complete tasks compared to the other baselines.

Table IV also includes a comparison of success rates and median completion times between an expert and new users utilizing OPEN TEACH for all tasks. On average, new users exhibit a success rate that is 76% of the expert's and take $2.25\times$ longer to complete a task. Additional details regarding individual user performances are provided in Appendix C. Intriguingly, some new users, despite their unfamiliarity with the framework, achieve comparable or superior performance to the experts in certain tasks. This observation highlights two factors: (1) the inherent variation in abilities among individuals, and (2) while our system is intuitive for new users, prolonged training leads to substantial improvement in their performance, with the potential for further enhancement with continued practice.

## VI. LIMITATIONS AND DISCUSSION

In this work, we introduce OPEN TEACH, an open-source unified framework designed to facilitate low-latency, high-frequency robot teleoperation. This versatile framework is tailored to accommodate diverse tasks and is compatible with a range of robot morphologies. However, we recognize a few limitations in this work: (*a*) OPEN TEACH relies on the accuracy of the in-built hand pose detection in the VR headset. Inaccuracies, particularly when fingers are occluded from view, can diminish the quality of hand tracking, posing

challenges to teleoperation. (*b*) In specific instances, the pose detector on the Oculus board may misconstrue finger positions, leading to difficulties in executing gestures like gripper closing, which relies on precise pinches between fingers. Addressing these challenges through future research on hand pose detection and tracking holds the potential to enhance the ease and intuitiveness of teleoperation using VR headsets.

## VII. Invitation for Contributions to Open Teach

We firmly advocate for the advancement of robotics through the unrestricted accessibility of research projects to the broader community. In line with this principle, we will open-source every component of this research endeavor, including the VR application, the human-robot interface, and the robot controllers. Further, we encourage fellow researchers in establishing the infrastructure for their own robotic setups, and we invite inquiries for assistance. Finally, we view our work as a step towards achieving more accessible and affordable robot teleoperation. Recognizing that there is ample room for enhancement, we enthusiastically welcome contributions to our repositories and would be happy to share such contributions with the world with proper attribution given to the contributors.

## References

[1] Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning. In *ICML*, 2004.

[2] Dafni Antotsiou, Guillermo Garcia-Hernando, and Tae-Kyun Kim. Task-oriented hand motion retargeting for dexterous manipulation imitation. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018.

[3] Sridhar Pandian Arunachalam, Sneha Silwal, Ben Evans, and Lerrel Pinto. Dexterous imitation made easy: A learning-based framework for efficient dexterous manipulation. *arXiv preprint arXiv:2203.13251*, 2022.

[4] Sridhar Pandian Arunachalam, Irmak Güzey, Soumith Chintala, and Lerrel Pinto. Holo-dex: Teaching dexterity with immersive mixed reality. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5962–5969. IEEE, 2023.

[5] Homanga Bharadhwaj, Jay Vakil, Mohit Sharma, Abhinav Gupta, Shubham Tulsiani, and Vikash Kumar. Roboagent: Generalization and efficiency in robot manipulation via semantic augmentations and action chunking. *arXiv preprint arXiv:2309.01918*, 2023.

[6] Aude G Billard, Sylvain Calinon, and Florent Guenter. Discriminative and adaptive imitation in uni-manual and bi-manual tasks. *Robotics and Autonomous Systems*, 54 (5):370–384, 2006.

[7] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.

[8] Manuel Caeiro-Rodríguez, Iván Otero-González, Fernando A. Mikic-Fonte, and Martín Llamas-Nistal. A systematic review of commercial smart gloves: Current status and applications. *Sensors*, 2021. ISSN 1424-8220. doi: 10.3390/s21082667.

[9] Xuxin Cheng, Kexin Shi, Ananye Agarwal, and Deepak Pathak. Extreme parkour with legged robots. *arXiv preprint arXiv:2309.14341*, 2023.

[10] Cheng Chi, Benjamin Burchfiel, Eric Cousineau, Siyuan Feng, and Shuran Song. Iterative residual policy: for goal-conditioned dynamic manipulation of deformable objects. *arXiv preprint arXiv:2203.00663*, 2022.

[11] Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *arXiv preprint arXiv:2303.04137*, 2023.

[12] Samuel Cohen, Brandon Amos, Marc Peter Deisenroth, Mikael Henaff, Eugene Vinitsky, and Denis Yarats. Imitation learning from pixel observations for continuous control. 2021.

[13] Zichen Jeff Cui, Yibin Wang, Nur Muhammad, Lerrel Pinto, et al. From play to policy: Conditional behavior generation from uncurated robot data. *arXiv preprint arXiv:2210.10047*, 2022.

[14] Hongjie Fang, Hao-Shu Fang, Yiming Wang, Jieji Ren, Jingjing Chen, Ruo Zhang, Weiming Wang, and Cewu Lu. Low-cost exoskeletons for learning whole-arm manipulation in the wild. *arXiv preprint arXiv:2309.14975*, 2023.

[15] Zipeng Fu, Tony Z Zhao, and Chelsea Finn. Mobile aloha: Learning bimanual mobile manipulation with low-cost whole-body teleoperation. *arXiv preprint arXiv:2401.02117*, 2024.

[16] Dhiraj Gandhi, Lerrel Pinto, and Abhinav Gupta. Learning to fly by crashing. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3948–3955. IEEE, 2017.

[17] Siddhant Gangapurwala, Mathieu Geisert, Romeo Orsolino, Maurice Fallon, and Ioannis Havoutis. Rloc: Terrain-aware legged locomotion using reinforcement learning and optimal control. *IEEE Transactions on Robotics*, 2022.

[18] Abraham George, Alison Bartsch, and Amir Barati Farimani. Openvr: Teleoperation for manipulation. *arXiv preprint arXiv:2305.09765*, 2023.

[19] Zaid Gharaybeh, Howard Chizeck, and Andrew Stewart. Telerobotic control in virtual reality. In *OCEANS 2019 MTS/IEEE SEATTLE*, pages 1–8, 2019. doi: 10.23919/ OCEANS40490.2019.8962616.

[20] Irmak Guzey, Yinlong Dai, Ben Evans, Soumith Chintala, and Lerrel Pinto. See to touch: Learning tac-

tile dexterity through visual incentives. *arXiv preprint arXiv:2309.12300*, 2023.

[21] Siddhant Haldar, Vaibhav Mathur, Denis Yarats, and Lerrel Pinto. Watch and match: Supercharging imitation with regularized optimal transport. *arXiv preprint arXiv:2206.15469*, 2022.

[22] Siddhant Haldar, Jyothish Pari, Anant Rai, and Lerrel Pinto. Teach a robot to fish: Versatile imitation from one minute of demonstrations. *arXiv preprint arXiv:2303.01497*, 2023.

[23] Shangchen Han, Beibei Liu, Randi Cabezas, Christopher D. Twigg, Peizhao Zhang, Jeff Petkau, Tsz-Ho Yu, Chun-Jung Tai, Muzaffer Akbay, Zheng Wang, Asaf Nitzan, Gang Dong, Yuting Ye, Lingling Tao, Chengde Wan, and Robert Wang. Megatrack: Monochrome egocentric articulated hand-tracking for virtual reality. 2020.

[24] Ankur Handa, Karl Van Wyk, Wei Yang, Jacky Liang, Yu-Wei Chao, Qian Wan, Stan Birchfield, Nathan Ratliff, and Dieter Fox. Dexpilot: Vision-based teleoperation of dexterous robotic hand-arm system. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9164–9170, 2020. doi: 10.1109/ICRA40945.2020.9197124.

[25] Thomas Hulin, Katharina Hertkorn, Philipp Kremer, Simon Schätzle, Jordi Artigas, Mikel Sagardia, Franziska Zacharias, and Carsten Preusche. The dlr bimanual haptic device with optimized workspace. In *2011 IEEE International Conference on Robotics and Automation*, pages 3441–3442. IEEE, 2011.

[26] Jemin Hwangbo, Inkyu Sa, Roland Siegwart, and Marco Hutter. Control of a quadrotor with reinforcement learning. *IEEE Robotics and Automation Letters*, 2(4):2096–2103, 2017. doi: 10.1109/LRA.2017.2720851.

[27] Benjamin G Katz. *A low cost modular actuator for dynamic robots*. PhD thesis, Massachusetts Institute of Technology, 2018.

[28] Sung-Kyun Kim, Seokmin Hong, and Doik Kim. A walking motion imitation framework of a humanoid robot by human walking recognition from imu motion data. In *2009 9th IEEE-RAS International Conference on Humanoid Robots*, pages 343–348. IEEE, 2009.

[29] Vikash Kumar and Emanuel Todorov. Mujoco haptix: A virtual reality system for hand manipulation. In *2015 IEEE-RAS 15th International Conference on Humanoid Robots (Humanoids)*, pages 657–663, 2015. doi: 10.1109/HUMANOIDS.2015.7363441.

[30] Marco Laghi, Michele Maimeri, Mathieu Marchand, Clara Leparoux, Manuel Catalano, Arash Ajoudani, and Antonio Bicchi. Shared-autonomy control for intuitive bimanual tele-manipulation. In *2018 IEEE-RAS 18th International Conference on Humanoid Robots (Humanoids)*, pages 1–9. IEEE, 2018.

[31] Shuang Li, Xiaojian Ma, Hongzhuo Liang, Michael Görner, Philipp Ruppel, Bin Fang, Fuchun Sun, and Jianwei Zhang. Vision-based teleoperation of shadow dexterous hand using end-to-end deep neural network.

In *2019 International Conference on Robotics and Automation (ICRA)*, pages 416–422. IEEE, 2019.

[32] Shuang Li, Jiaxi Jiang, Philipp Ruppel, Hongzhuo Liang, Xiaojian Ma, Norman Hendrich, Fuchun Sun, and Jianwei Zhang. A mobile robot hand-arm teleoperation system by vision and imu. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 10900–10906. IEEE, 2020.

[33] Shuang Li, Norman Hendrich, Hongzhuo Liang, Philipp Ruppel, Changshui Zhang, and Jianwei Zhang. A dexterous hand-arm teleoperation system based on hand pose estimation and active vision. *IEEE Transactions on Cybernetics*, 2022.

[34] Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu, Yuke Zhu, and Peter Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. *arXiv preprint arXiv:2306.03310*, 2023.

[35] Shaowei Liu, Hanwen Jiang, Jiarui Xu, Sifei Liu, and Xiaolong Wang. Semi-supervised 3d hand-object poses estimation with interactions in time. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14687–14697, 2021.

[36] Yuntao Ma, Farbod Farshidian, Takahiro Miki, Joonho Lee, and Marco Hutter. Combining learning-based locomotion policy with model-based manipulation for legged mobile manipulators. *IEEE Robotics and Automation Letters*, 7(2):2377–2384, 2022. doi: 10.1109/LRA.2022.3143567.

[37] Ajay Mandlekar, Yuke Zhu, Animesh Garg, Jonathan Booher, Max Spero, Albert Tung, Julian Gao, John Emmons, Anchit Gupta, Emre Orbay, et al. Roboturk: A crowdsourcing platform for robotic skill learning through imitation. In *Conference on Robot Learning*, pages 879–893. PMLR, 2018.

[38] Cassie Meeker, Maximilian Haas-Heger, and Matei Ciocarlie. A continuous teleoperation subspace with empirical and algorithmic mapping algorithms for nonanthropomorphic hands. *IEEE Transactions on Automation Science and Engineering*, 19(1):373–386, 2020.

[39] Malte Mosbach, Kara Moraw, and Sven Behnke. Accelerating interactive human-like manipulation learning with gpu-based simulation and high-quality demonstrations. In *2022 IEEE-RAS 21st International Conference on Humanoid Robots (Humanoids)*, pages 435–441. IEEE, 2022.

[40] Ashvin Nair, Abhishek Gupta, Murtaza Dalal, and Sergey Levine. Awac: Accelerating online reinforcement learning with offline datasets. *arXiv preprint arXiv:2006.09359*, 2020.

[41] Andrew Y Ng, Stuart J Russell, et al. Algorithms for inverse reinforcement learning. In *ICML*, 2000.

[42] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Charles Xu, Jianlan Luo, Tobias Kreiman, You Liang Tan, Dorsa Sadigh, Chelsea Finn, and Sergey Levine. Octo: An open-source generalist robot policy.

https://octo-models.github.io, 2023.

[43] Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anikait Singh, Anthony Brohan, et al. Open x-embodiment: Robotic learning datasets and rt-x models. *arXiv preprint arXiv:2310.08864*, 2023.

[44] Jyothish Pari, Nur Muhammad Shafiullah, Sridhar Pandian Arunachalam, and Lerrel Pinto. The surprising effectiveness of representation learning for visual imitation, 2021.

[45] Dean A. Pomerleau. Alvinn: An autonomous land vehicle in a neural network. In D. Touretzky, editor, *NeurIPS*, volume 1. Morgan-Kaufmann, 1988.

[46] Yuzhe Qin, Hao Su, and Xiaolong Wang. From one hand to multiple hands: Imitation learning for dexterous manipulation from single-camera teleoperation. *arXiv preprint arXiv:2204.12490*, 2022.

[47] Yuzhe Qin, Wei Yang, Binghao Huang, Karl Van Wyk, Hao Su, Xiaolong Wang, Yu-Wei Chao, and Dietor Fox. Anyteleop: A general vision-based dexterous robot arm-hand teleoperation system. *arXiv preprint arXiv:2307.04577*, 2023.

[48] Ilija Radosavovic, Tete Xiao, Stephen James, Pieter Abbeel, Jitendra Malik, and Trevor Darrell. Real-world robot learning with masked visual pre-training, 2022. URL https://arxiv.org/abs/2210.03109.

[49] Moritz Reuss, Maximilian Li, Xiaogang Jia, and Rudolf Lioutikov. Goal-conditioned imitation learning using score-based diffusion policies. *arXiv preprint arXiv:2304.02532*, 2023.

[50] Douglas A Reynolds et al. Gaussian mixture models. *Encyclopedia of biometrics*, 741(659-663), 2009.

[51] Max Schwarz, Christian Lenz, Andre Rochow, Michael Schreiber, and Sven Behnke. Nimbro avatar: Interactive immersive telepresence with force-feedback telemanipulation. in 2021 ieee. In *RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5312–5319.

[52] Nur Muhammad Shafiullah, Zichen Cui, Ariuntuya Arty Altanzaya, and Lerrel Pinto. Behavior transformers: Cloning $k$ modes with one stone. *Advances in neural information processing systems*, 35:22955–22968, 2022.

[53] Nur Muhammad Mahi Shafiullah, Anant Rai, Haritheja Etukuru, Yiqian Liu, Ishan Misra, Soumith Chintala, and Lerrel Pinto. On bringing robots home. *arXiv preprint arXiv:2311.16098*, 2023.

[54] Dhruv Shah, Ajay Sridhar, Nitish Dashora, Kyle Stachowicz, Kevin Black, Noriaki Hirose, and Sergey Levine. ViNT: A foundation model for visual navigation. In *7th Annual Conference on Robot Learning*, 2023. URL https://arxiv.org/abs/2306.14846.

[55] Neo Ee Sian, Kazuhito Yokoi, Shuuji Kajita, Fumio Kanehiro, and Kazuo Tanie. Whole body teleoperation of a humanoid robot development of a simple master device using joysticks. *Journal of the Robotics Society of Japan*, 22(4):519–527, 2004.

[56] Aravind Sivakumar, Kenneth Shaw, and Deepak Pathak. Robotic telekinesis: Learning a robotic hand imitator by watching humans on youtube, 2022.

[57] Laura Smith, Ilya Kostrikov, and Sergey Levine. A walk in the park: Learning to walk in 20 minutes with model-free reinforcement learning. *arXiv preprint arXiv:2208.07860*, 2022.

[58] Shuran Song, Andy Zeng, Johnny Lee, and Thomas Funkhouser. Grasping in the wild: Learning 6dof closed-loop grasping from low-cost demonstrations. *RA-L*, 2020.

[59] Stanford Artificial Intelligence Laboratory et al. Robotic operating system. URL https://www.ros.org.

[60] Chen Wang, Linxi Fan, Jiankai Sun, Ruohan Zhang, Li Fei-Fei, Danfei Xu, Yuke Zhu, and Anima Anandkumar. Mimicplay: Long-horizon imitation learning by watching human play. *arXiv preprint arXiv:2302.12422*, 2023.

[61] Philipp Wu, Yide Shentu, Zhongke Yi, Xingyu Lin, and Pieter Abbeel. Gello: A general, low-cost, and intuitive teleoperation framework for robot manipulators. *arXiv preprint arXiv:2309.13037*, 2023.

[62] Yuqiang Wu, Pietro Balatti, Marta Lorenzini, Fei Zhao, Wansoo Kim, and Arash Ajoudani. A teleoperation interface for loco-manipulation control of mobile collaborative robotic assistant. *IEEE Robotics and Automation Letters*, 4(4):3593–3600, 2019.

[63] xArm Developer. xarm python sdk. https://github.com/xArm-Developer/xArm-Python-SDK.

[64] Fan Zhang, Valentin Bazarevsky, Andrey Vakunov, Andrei Tkachenka, George Sung, Chuo-Ling Chang, and Matthias Grundmann. Mediapipe hands: On-device real-time hand tracking, 2020.

[65] Tianhao Zhang, Gregory Kahn, Sergey Levine, and Pieter Abbeel. Learning deep control policies for autonomous aerial vehicles with mpc-guided policy search. In *2016 IEEE international conference on robotics and automation (ICRA)*, pages 528–535. IEEE, 2016.

[66] Tianhao Zhang, Zoe McCarthy, Owen Jow, Dennis Lee, Xi Chen, Ken Goldberg, and Pieter Abbeel. Deep imitation learning for complex manipulation tasks from virtual reality teleoperation. In *ICRA*, 2018.

[67] Tony Z Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual manipulation with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.

[68] Wenping Zhao, Jinxiang Chai, and Ying-Qing Xu. Combining marker-based mocap and rgb-d camera for acquiring high-fidelity hand motion data. In *Proceedings of the ACM SIGGRAPH/eurographics symposium on computer animation*, pages 33–42, 2012.

[69] Yifeng Zhu, Abhishek Joshi, Peter Stone, and Yuke Zhu. Viola: Imitation learning for vision-based manipulation with object proposal priors. *arXiv preprint arXiv:2210.11339*, 2022. doi: 10.48550/arXiv.2210.11339.

## APPENDIX

### A. Framework details

*1) Structure of the framework:* We use ZeroMQ for networking between nodes. The OPEN TEACH framework is divided into two parts - *teleoperation* and *data collection*.

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

1) **Detector:** Receives the hand keypoints from the Meta Quest 3 and publishes them to ZMQ sockets.
2) **Keypoint Transformer:** Subscribes the keypoints published by the detector and maps them to the robot pose.
3) **Operator:** Receives the robot pose from the keypoint transformer and the current robot state from the controller. The operator computes the robot's actions which are published to a ZMQ socket.
4) **Controller:** Subscribes an action from the operator and takes an action in the real or simulated environment. After taking the action, the controller publishes the current states of the environment for use by the operator.
5) **Visualizer:** Subscribes the RGB images from the camera process (or the environment in case of simulations) and puts it on the screen inside the VR app for visualization during teleoperation.

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

*2) Thumb Retargeting for Robot Hand:* Section IV-C provides details about the design of the OPEN TEACH wrapper for the robot hand. To recap, given the individual joint angles in the teacher's hand from the VR headset, the joint angles for the robot hand can be computed by directly commanding the robot's joints to the corresponding angles. This works well for all fingers except the thumb. Holo-Dex[4] deals with this by mapping the spatial coordinate of the teacher's thumb tip to that of the robot hand. Then an inverse kinematics solver is used to compute the joint angles of the thumb. In this case, the retargeting of the thumb is done in 2D space. These bounds, depicted in Fig. 5(a), define the thumb's reach limits. During retargeting, the thumb tip's zone on the 2D palm plane is detected, and a perspective transform from the human hand to the robot hand is applied, aligning the human thumb tip with the robot thumb tip on the 2D plane. However, using three separate bounds introduces jitters when the thumb tip transitions between zones and results in stagnancy when outside the bounds. Further, in Holo-Dex, the height of the robot thumb tip is fixed, allowing it to only move along the 2D space.

To address these challenges, OPEN TEACH employs a single, large zone spanning the entire thumb's workspace in 2D space(refer to Fig. 5(b)). When the thumb is within bounds, a perspective transformation retargets the human thumb tip to the robot thumb tip. In cases where the thumb goes out of bounds, the closest point within the bound is estimated and used for retargeting, avoiding stagnation. Additionally, instead of a fixed height, the thumb is allowed to move perpendicular to the 2D surface along the palm, mapping the height of the human thumb tip to the robot thumb tip based on maximum and minimum height bounds. This approach ensures smoother thumb motion and enables the performance of more complex tasks compared to Holo-Dex [4].

TABLE V: 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 |

### B. Task Details

*1) Demo Collections times:* Table V provides the average times required to collect a demonstration for 16 tasks across 3 real-world setups (Franka-Allegro, Kinova-Allegro, Bimanual) and 2 simulated environments(Allegro sim, LIBERO sim).

*2) Task Descriptions:* Fig. 6, Fig. 7, Fig. 8, Fig. 9, Fig. 10, and Fig. 11 provide rollouts of all the tasks performed both in the real world and in simulated environments. Each task rollout is labeled with the name of the task and a task description.

### C. User Study

Following up from Section V-F, we provide the success rate and average completion times for all 15 users for each task performed in Table VI and Table VII respectively. Each user roughly performed 3 tasks on average, with 5 trials for each task. As mentioned in Section V-F, since the Holo-Dex [4] and AnyTeleop [47] baselines lack open-source code for arm retargeting, we were unable to evaluate them on tasks involving arm movements. We observe a wide range of differences in success rates and average completion times demonstrating the inherent variations across users.

(a) Holo-Dex     (b) Open Teach

Fig. 5: Thumb retargeting difference

TABLE VI: 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 VII: 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 | - | - |

**Slice Cucumber:** Stabilize the cucumber with one arm and use a knife for cutting a slice with the other arm.

**Make Sandwich:** Make a sandwich by sequentially adding the ingredients placed on a table.

**Insert USB:** Insert a USB charging cable into the adapter.

**Toast Bread:** Place two slices of bread in the toaster oven.

**Sweep Dirt Off the Table:** Sweep dirt off the table using a brush and a dustpan.

**Unfold Cloth:** Unfold the cloth and place it on the table.

**Uncap Marker:** Lift the marker with one arm and uncap it with the other.

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

**Bimanual**

**Handover can:** Pick up the can with one arm and hand it over to the other arm.

**Pour Sprinkles:** Pour sprinkles into a bowl.

**Franka-Allegro**

**Pack and Lift a Basket:** Place the honey bottle and soda can inside the basket, lift the handles and pick up the basket.

**Put Muffin to the Oven:** Place a muffin inside the oven and close the door.

**Undo latch and open box:** Undo the latch and open the box.

**Slot battery:** Slot a battery into the battery holder.

**Shell Game:** Play the shell game with 2 cups and a ball.

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

**Franka-Allegro**

**Open box:** Lift the lid of the box.

**Grasp object and twist:** Grasp the lemon on the table and rotate the hand to face upwards.

**Grasp Sponge:** Grasp sponge placed on the table.

**Pick Up Tea Sachet:** Slide the tea sachet and pick it up.

**Kinova-Allegro**

**Stack Blocks:** Stack blocks on top of each other.

**Open Ketchup Bottle:** Open the cap of a ketchup bottle.

**Open Drawer:** Slide the drawer and open it.

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

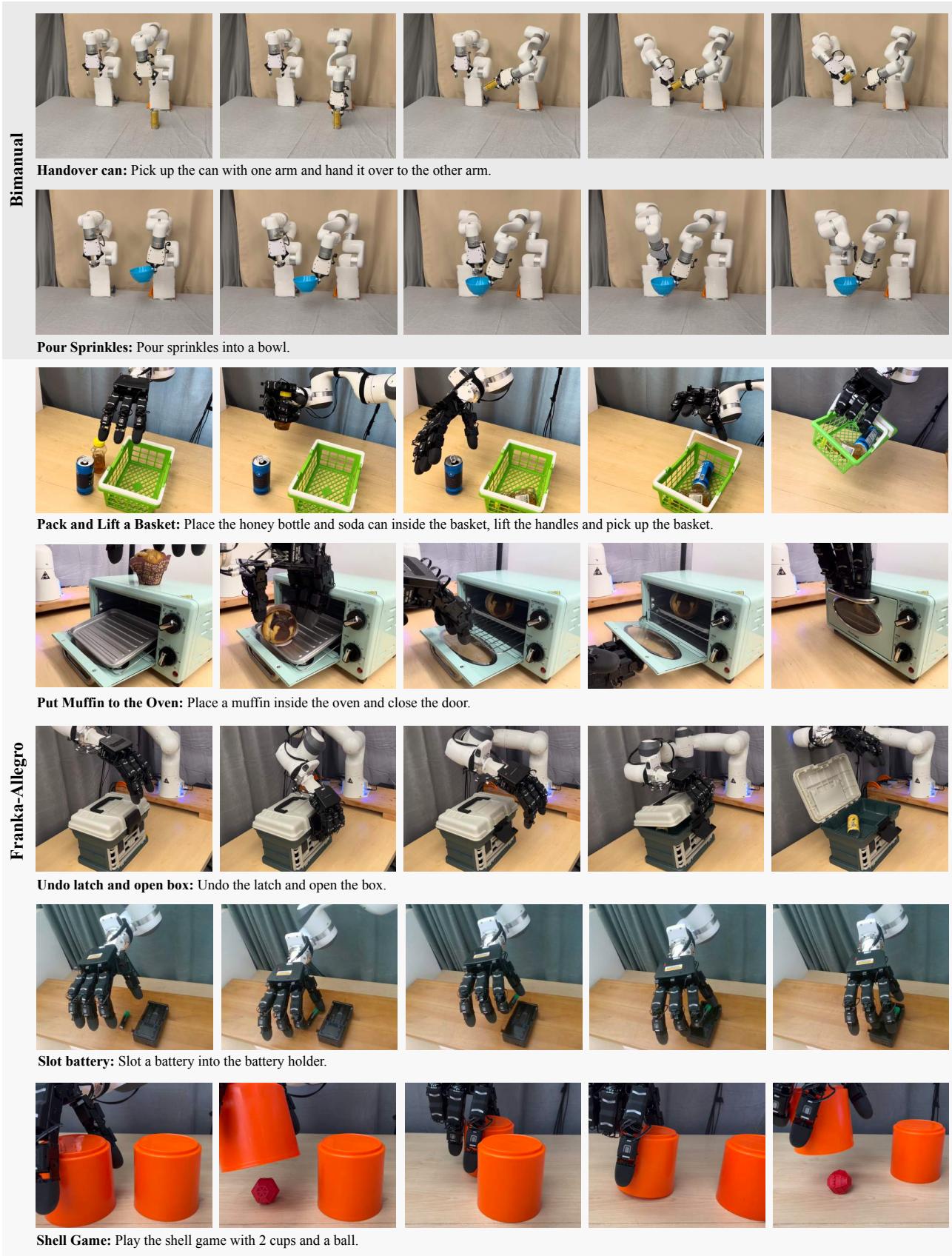**Iron Towel:** Move the hand to the iron and iron the towel.

**Hit Hammer:** Hit a nail with the hammer.

**Laptop Opening:** Open the lid of a laptop. The passthrough of the VR app is being streamed on the laptop's screen.

**Unfold Towel:** Unfold the tower placed on the table.

**Undo Latch and Open Box:** Undo the latch and open the box.

**Write Alphabet:** Write the alphabet "A" on a paper placed on the table.

**Flip Sponge:** Flip the sponge on the table.

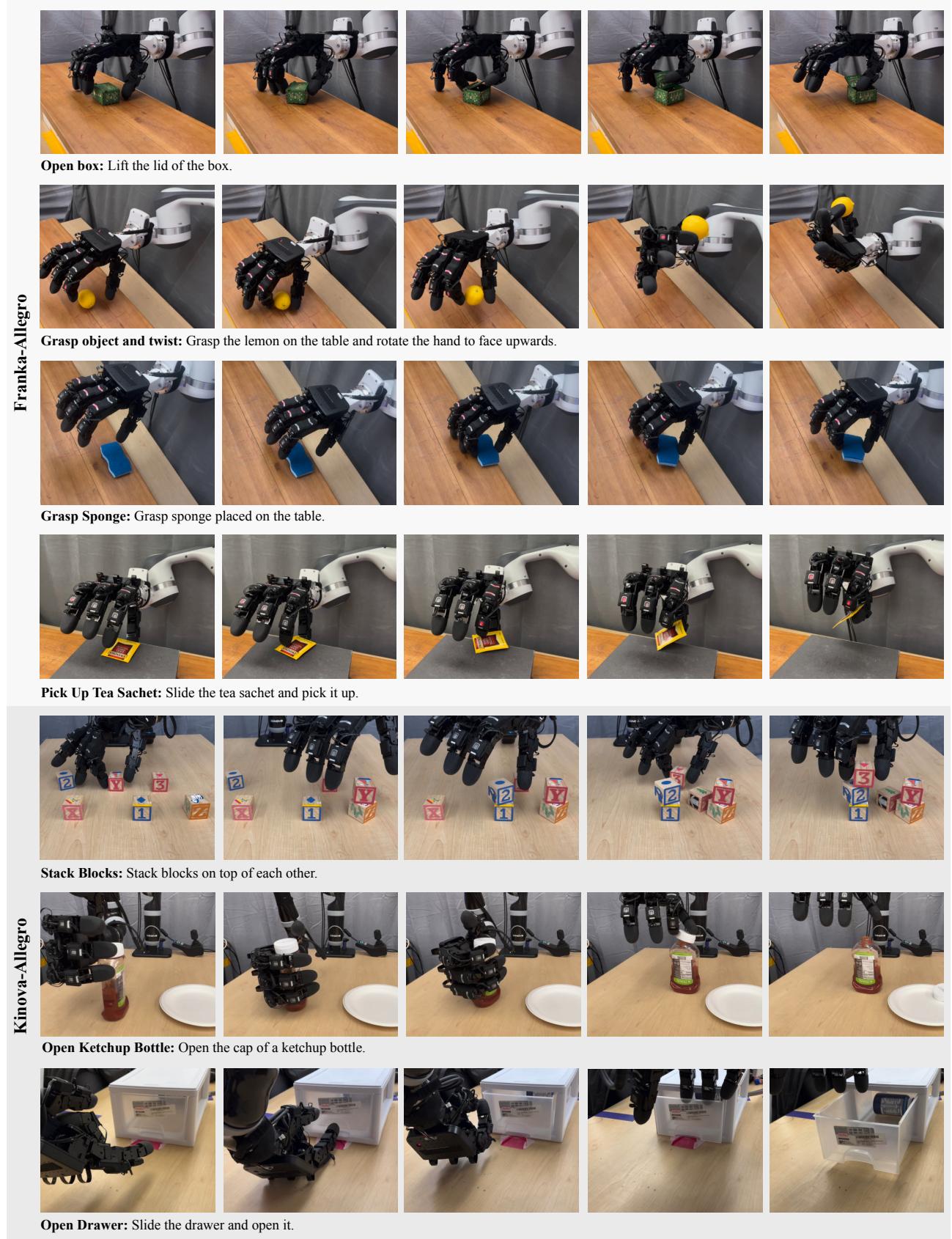Fig. 9: Real world task rollouts demonstrating the ability of OPEN TEACH to perform intricate, long-horizon tasks.

**Allegro Sim**
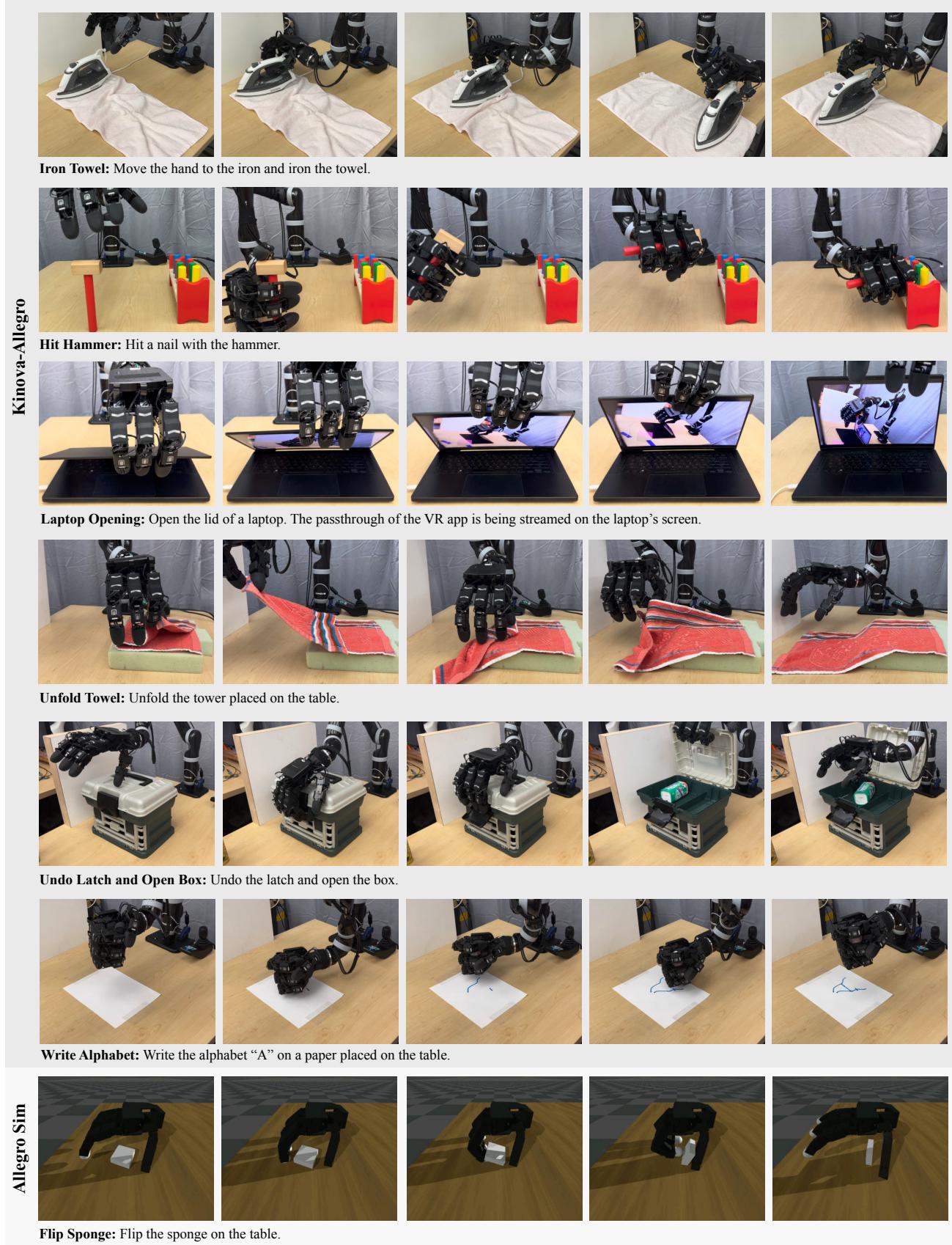
**Pinch Grasp:** Grasp the block with a pinch and lift it.

**Flip Cube:** Flip the cube in-hand.

**LIBERO Sim**

**Turn on Stove:** Turn on the stove.

**Place white mug on left plate:** Pick up the white mug and place it on the left white plate.

**Put pan on stove and turn it on:** Pick up the pan and put it on the stove. Then turn on the stove.

**Organize basket:** Pick up everything placed on the table and put it in the basket.

**Close top drawer of cabinet:** Approach the top drawer of the cabinet and put it to close it.

Fig. 10: Real world task rollouts demonstrating the ability of OPEN TEACH to perform intricate, long-horizon tasks.

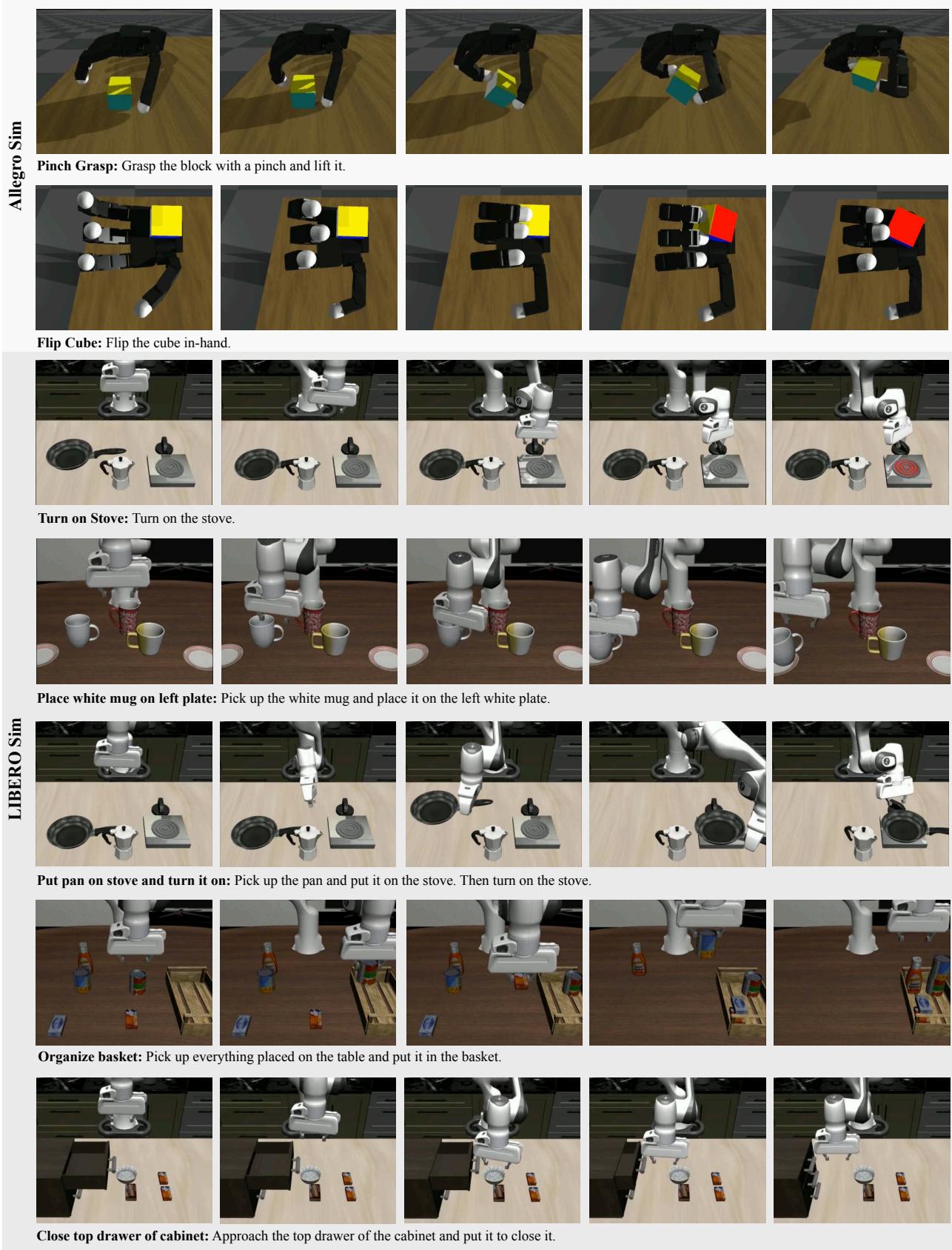**Open Drawer:** Move close to the drawer and open it.

**Open Door:** Move close to the cabinet door and open it.

**Put plastic bag in garbage can:** Pick up the plastic bag from the counter and put it in the garbage can.
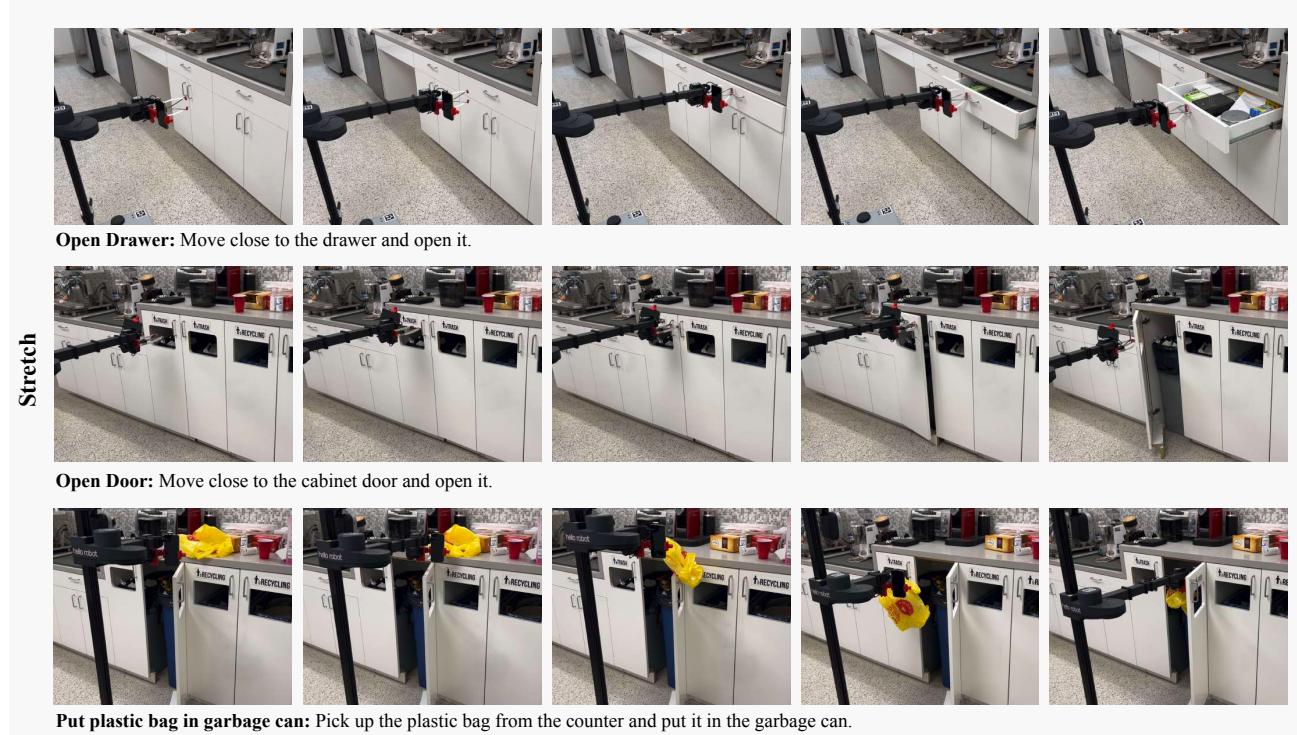
Fig. 11: Real world task rollouts demonstrating the ability of OPEN TEACH to perform intricate, long-horizon tasks.