

TeleMoMa: A Modular and Versatile Teleoperation System for Mobile Manipulation

Author Names Omitted for Anonymous Review.

Abstract—A critical bottleneck limiting imitation learning in robotics is the lack of data. This problem is more severe in mobile manipulation, where collecting demonstrations is harder than in stationary manipulation due to the lack of available and easy-to-use teleoperation interfaces. In this work, we demonstrate TeleMoMa, a general and modular interface for whole-body teleoperation of mobile manipulators. TeleMoMa unifies multiple human interfaces including RGB and depth cameras, virtual reality controllers, keyboard, joysticks, etc., and any combination thereof. In its more accessible version, TeleMoMa works using simply vision (e.g., an RGB-D camera), lowering the entry bar for humans to provide mobile manipulation demonstrations. We demonstrate the versatility of TeleMoMa by teleoperating several existing mobile manipulators — PAL Tiago++, Toyota HSR, and Fetch — in simulation and the real world. We demonstrate the quality of the demonstrations collected with TeleMoMa by training imitation learning policies for mobile manipulation tasks involving synchronized whole-body motion. Finally, we also show that TeleMoMa’s teleoperation channel enables teleoperation *on site*, looking at the robot, or *remote*, sending commands and observations through a computer network, and perform user studies to evaluate how easy it is for novice users to learn to collect demonstrations with different combinations of human interfaces enabled by our system. We hope TeleMoMa becomes a helpful tool for the community enabling researchers to collect whole-body mobile manipulation demonstrations.

I. INTRODUCTION

A core goal of robotics is to build generalist robots capable of operating alongside humans in their environment. To this end, learning from human-collected robot demonstrations has shown promise in endowing robots with the capabilities to solve complex tasks [4, 42, 32], boosted recently by the advent of foundation models capable of learning from large amounts of data [5, 14]. While these models demonstrate an impressive semantic understanding of the tasks [51, 6, 9, 12, 8], these successes have been largely limited to stationary manipulation. However, a large fraction of the tasks that we would like generalist robots to perform require a combination of manipulation and mobility: e.g., sweeping the floor requires moving the broom with both hands and walking around to reach the dirty spots; covering a table with a tablecloth requires holding the tablecloth and pulling it over the table while simultaneously moving to reach all edges.

One of the reasons why stationary manipulation has enjoyed the benefits of large models, while mobile manipulation has not, is due to the availability of large datasets of human-collected demonstrations [9, 54]. They were obtained due to the multiple existing and easy-to-use teleoperation frameworks for stationary manipulators [30, 64, 59, 39, 10]. For mobile manipulation, however, the existing stationary manipulation



Fig. 1: **TeleMoMa**: a modular and versatile mobile manipulation teleoperation system. *Left and Middle* Demonstrators performing a bimanual sweeping task with the vision-only, virtual reality (VR), and a combination of vision+VR interfaces. TeleMoMa enables multiple human interfaces and their combination. *Middle and Right* Tiago (real), HSR (real), and Fetch (simulation), three of the robot platforms that we demonstrate teleoperated for different mobile manipulation tasks with TeleMoMa, demonstrating its versatility.

teleoperation systems are not sufficient, due to the additional degrees of freedom that the user has to control including mobility and possibly multiple arms.

Several teleoperation frameworks for mobile manipulation have been proposed in the past, with different capabilities and limitations. They either enable accurate control with specific (and often expensive) hardware like motion capture systems [3, 50, 47] or puppeteering interfaces [16, 38], or achieve scalability by overloading simple and available devices that work for stationary manipulators such as gamepads [10], virtual reality controllers [46], or mobile phones [52, 58], limiting the expressiveness of the demonstrations. Teleoperation based solely on vision [63, 39, 21] promises an available and accessible interface at the cost of accuracy and dexterity. Each device alone presents a tradeoff between accuracy and availability, versatility and expressiveness, and as a result, no single device enables scalable, expressive teleoperation for all mobile manipulators.

Inspired by the complementary capabilities of several of the human interfaces for teleoperation, we introduce TeleMoMa (**Teleoperation for Mobile Manipulation**), a novel teleoperation framework focusing on modularity and versatility. TeleMoMa enables users to teleoperate different mobile manipulators with a variety of human interfaces or combi-

nations thereof, in simulation or the real world, providing users the means to select the combination that best fits their teleoperation needs. TeleMoMa offers the lowest entry point for researchers: it enables whole-body teleoperation with just a depth camera. But its modularity enables us to use other interfaces such as virtual reality (VR) controllers, a keyboard, a 3D mouse, a mobile phone, or their combination with vision, overcoming the limitations of each individual input modality. Such a system enables researchers to collect demonstrations of whole-body mobile manipulation tasks at scale for virtually any robot and hardware interface available. We demonstrate that TeleMoMa allows researchers to teleoperate different mobile manipulation platforms out-of-box such as PAL Tiago [36], Toyota HSR [60] and Zebra Fetch [57]. TeleMoMa extends also to simulation, which we demonstrate by integrating it with OmniGibson [27] and the BEHAVIOR-1K benchmark.

In our experiments, we evaluate both TeleMoMa’s usability and its suitability for data collection for imitation learning. We conducted a user study to evaluate the benefits of modularity in TeleMoMa and accessibility of TeleMoMa-enabled interfaces for novice users. Our results indicate that a hybrid vision-VR interface is an efficient and natural mode of teleoperation, and that novice users are quickly able to learn to use it. We also successfully trained several imitation learning policies on the data collected using TeleMoMa and explored relevant questions for IL for mobile manipulation such as (a) *What inputs matter in imitation learning for mobile manipulators?*, and (b) *How do the policies scale as we increase the size of the training data?* We measured the role that different embodiments with different capabilities and the sim-real gap have in teleoperation performance, and demonstrate remote teleoperation of real robots with an analysis of robustness to the latency and delays in the communication.

In summary, TeleMoMa is a novel, modular, and versatile teleoperation framework for whole-body mobile manipulation that facilitates the integration of different human interfaces, robot platforms, and simulators. We hope that our contribution lowers the barrier of entry for researchers to collect demonstrations for imitation learning for mobile manipulation.

II. RELATED WORK

Teleoperation for General Robotics. Teleoperation is almost as old as the field of robotics itself [53], with early manipulators being controlled in kinematically identical master-slave systems [48] similar to the very recent Mobile ALOHA [16]. More recently, teleoperation has emerged as a critical means of data collection for imitation learning methods [4, 42], as the ability to quickly collect large scale robotic data has become paramount for training large capacity behavior models [31, 22, 10]. Many teleoperation modalities have been proposed to address these challenges, including kinesthetic teaching, joysticks, virtual reality, mobile phones, RGB cameras, exoskeletons, and motion capture.

Each modality has its benefits and shortcomings. Joysticks (e.g. the SpaceMouse) offer intuitive control of a robot’s end-

effector(s), but fail to enable joint control or navigation [43]. Virtual reality enables users to perform tasks from the robot’s perspective, but is limited by individual tolerance to motion sickness and does not naturally enable simultaneous locomotion and manipulation [62, 56, 40, 11, 19, 28]. Mobile phones offer scalable data collection, but provide a very limited interface, failing to naturally support joint control or base motion [30, 31]. RGB cameras have been explored as an accessible, scalable medium with limited mobility and range of motion [21, 39, 49]. Exoskeletons and master-slave devices enable dexterous control but are typically platform-specific and costly [13, 64, 16, 38], and do not naturally provide a way to coordinate base and arm motion. Motion capture similarly enables high-quality data collection, but is costly and difficult to scale [3, 47, 50]. Kinesthetic teaching was the predominant teleoperation paradigm for imitation learning for many years [4, 42, 24], but fails to enable more complicated bimanual or mobile manipulation tasks. Some works explore the combination of different modalities [55, 15] but fail to be sufficiently general and extensible. Thus, despite the plethora of available options, there remains a need for a teleoperation system capable of adapting to the needs of mobile manipulation in a scalable, accessible way.

Teleoperation for Mobile Manipulation. Recent successes in learning from large collections of human demonstrations has been limited to stationary manipulators [9, 51, 6] or simple mobile manipulation tasks like pick and place that do not require coordination between base and arm motion [7, 8]. This is in part due to the lack of accessible and intuitive ways to collect demonstrations for mobile robots. Recently, some methods have tried to address this using specialized hardware, such as motion capture systems [3, 50, 25, 47], exoskeletons [16, 61, 35, 13] and more sophisticated human-computer interfaces [45, 26]. On the other hand, several works borrow from successful teleoperation interfaces in stationary manipulation, using interfaces such as VR [17, 46, 37, 20, 18, 34, 23], kinesthetic teaching [61], visual motion tracking [63], keyboard and mouse [41] and mobile phones [58], by modifying them to enable the control of mobile manipulators. Although these interfaces are accessible, they lack the granularity necessary to coordinate all degrees of freedom of a mobile robot for a true mobile manipulation task.

We summarize the main features of TeleMoMa and contrast it with related systems in Table I. The criteria for each category is described further in Appendix A. We compare across two primary dimensions: the teleoperation modalities provided, and the robot capabilities enabled. TeleMoMa is the only teleoperation system to provide modularity and enable the flexible combination of multiple input modalities. Moreover, it is the only system capable of full whole-body motion (including torso control) that remains accessible and robot agnostic.

III. TELEMOMA SYSTEM

TeleMoMa is a teleoperation system for mobile manipulators — versatile for different robot morphologies and modular

TABLE I: Comparison of Existing Mobile Manipulation Teleoperation Systems

	Teleoperation Support			Robot Support					
	Cost / Accessibility	Modular	Modality	Bimanual	Height Control	Whole-Body Teleop	Robot Agnostic	Action Space	Domain
Arduengo et al. [3]	\$\$	×	Mocap	×	✓	✓	✓	EE Pose / Base Vel.	Real
MoMaRT [58]	\$	×	Phone	×	×	×	✓	EE Pose / Base Vel.	Sim
MOMA-Force [61]	\$\$\$	×	Kinesthetic	×	×	×	✓	EE Pose and Wrench	Real
SATYRR [38]	\$\$\$	×	Puppeteer	✓	×	✓	×	Joint Pos. / Base Vel.	Real
TRILL [46]	\$	×	VR	✓	×	✓	✓	EE Poses / Gait	Sim&Real
Mobile ALOHA [16]	\$\$\$	×	Puppeteer	✓	×	✓	×	Joint Pos. / Base Vel.	Real
TeleMoMa	\$	✓	*	✓	✓	✓	✓	EE Poses / Base Vel. / Joint Pos.	Sim&Real

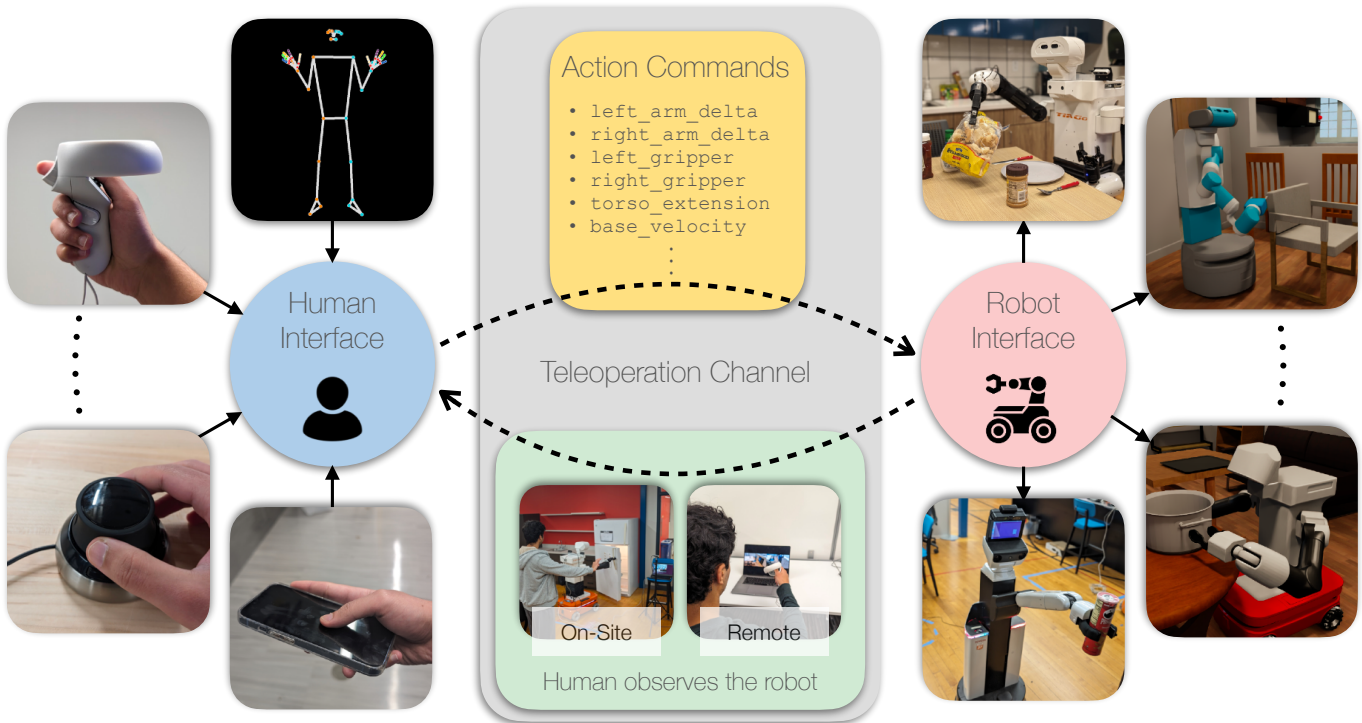


Fig. 2: **TeleMoMa System**. TeleMoMa consists of three components: the *Human Interface* acquires commands from the human using different input devices; the *Teleoperation Channel* defines the action command structure between the human and the robot interfaces, and, possibly, closes the loop with observations from the robot; and the *Robot Interface* implements a robot-specific mapping of actions to low-level robot commands. This architecture enables modularity and versatility – combining multiple devices to achieve intuitive whole-body teleoperation for multiple tasks and robots.

in the human input devices. It is generally composed of a *Teleoperation Channel* that defines the communication between a *Human Interface* and a *Robot Interface* (Fig. 2). The *Human Interface* acquires human inputs across different teleoperation modes such as vision, VR, spacemouse, keyboard, and mobile phones, or their combinations, and maps them to a general mobile manipulation action command structure provided by the *Teleoperation Channel* that includes fields such as base, arm, gripper, and torso motion. Multiple input devices can be combined through our *Human Interface* to acquire the action commands in the best suited manner for a task. The *Teleoperation Channel* hands over the action commands to the *Robot Interface*, a robot-specific module that maps the actions to robot motor commands. While the specific implementation

of the robot interface relies on platform-dependent controllers, our requirements (controllers for the motion of the end-effectors, base, joints, ...) are general enough to enable the teleoperation of most existing platforms, in the real world and simulation. In the following, we provide additional information about the three components of TeleMoMa.

A. Human Interface

The *Human Interface* is responsible for processing the captured data from various teleoperation input devices and mapping them to a common action command structure. For each input device, the data is processed independently by a device-specific parser that maps the signals from the input modality (keyboard strokes, motion of a VR controller, location of

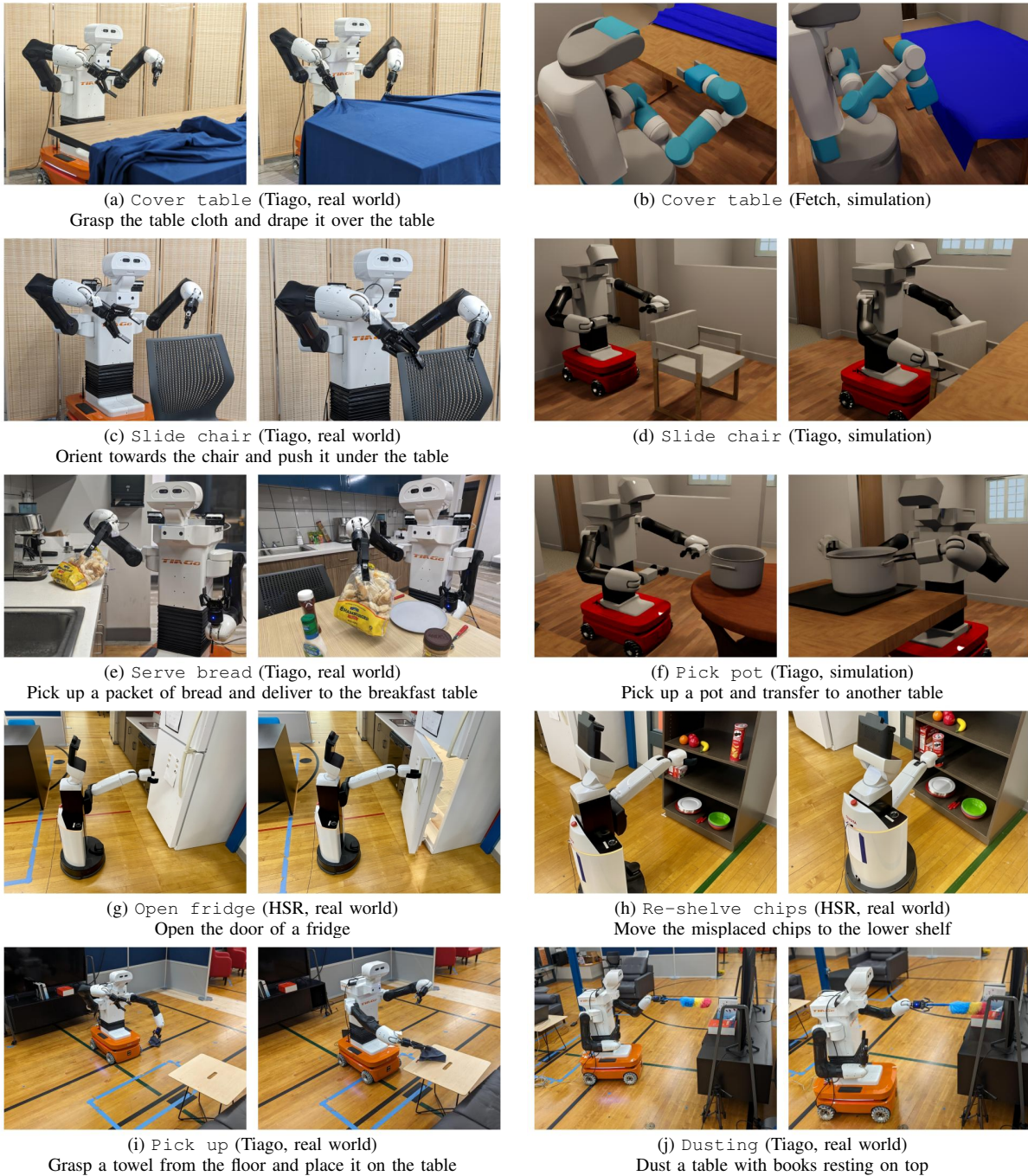


Fig. 3: **Tasks in our evaluation of TeleMoMa.** Shown above is the initial and goal state of each task.

human skeleton keypoints on an image, ...) into elements of the teleoperation channel's action command. TeleMoMa supports input modalities such as vision, keyboard, spacemouse, VR (Oculus Quest and HTC Vive) and mobile phones. In the following, we explain the vision and VR human interfaces from TeleMoMa in detail. We test these human interfaces (and their combination) extensively in our experiments because their combined capabilities strike a good balance between availability, generality, dexterity and accuracy. We defer the implementation details of other human interfaces such as

spacemouse, keyboard, and mobile phones to Appendix B.

1) *Vision-Based Human Interface:* TeleMoMa offers a unique vision-based pipeline for the whole-body teleoperation of a mobile manipulator using a single RGB-D camera. We use MediaPipe [29], a lightweight RGB-based model that executes in real-time for body pose and hand keypoint detection. Our proposed human interface uses the position and rotation of the hips to control the movement of the base of the mobile manipulator. Since the model only provides the relative depth of the keypoints to the center of the hip and not the absolute

depth, we use the depth channel of an RGB-D camera to obtain the absolute values. The hand keypoints are mapped to the end-effector of the robot based on the position and orientation of the palm with respect to the hip. We compute the per-frame relative pose displacement in Cartesian space of the hands and send them in the teleoperation channel’s action command as arm delta commands. Additionally, we use the distance between the center of the hips and ankles to command the robot height for robots with an actuated torso.

2) *Virtual Reality Controllers as Human Interface*: TeleMoMa supports Oculus Quest and HTC Vive virtual reality hardware devices as inputs to the VR human interface. The controllers are tracked with respect to the headset for Oculus and with respect to the lighthouse for HTC Vive. Similar to Seo et al. [46], the tracked hand poses in Cartesian space are used to command the end-effector in the task-space. As in the vision-based interface, we compute the per-frame relative pose displacement of hands and use them in the teleoperation channel’s action command. The joysticks integrated in the VR are used to command the velocities of the mobile base and also control the torso extension.

B. Teleoperation Channel

The *Teleoperation Channel* defines how the *Human Interface* communicates with the *Robot Interface*, and is the key to TeleMoMa’s generality and modularity. Specifically, the *Teleoperation Channel* defines an action command structure that serves as a bridge between the human and the robot and the way the active human interfaces populate the entries of this structure.

During deployment, users can specify what input modality they want to use to control each part of the robot’s embodiment including left and right arms and hands, torso, and base. The *Teleoperation Channel* automatically manages the action assignment based on the user specification, and consolidates the possible missing elements of the action commands due to differences in hardware frequency or network delays.

Finally, the *Teleoperation Channel* also defines the mechanism by which humans *close the loop* with the robot and observe the execution of the action commands, adapting those to achieve the mobile manipulation tasks. We consider two methods of observation: on-site and remote. When on-site, the human directly observes the robot executing the action commands. When remote, the *Teleoperation Channel* communicates the images from the onboard sensors of the robot to the human interface to be displayed for the human, enabling teleoperation from a different location. We evaluate both modes in our experiments (Sec. IV).

C. Robot Interface

The *Robot Interface* is a robot-specific module that maps the commands obtained from the *Human Interface* to the motor commands to the robot. In most of our experiments, those are torques at the joints of the robot. The specific controllers used to compute the torques are not part of the TeleMoMa system but they are necessary to map the action commands obtained

from the *Teleoperation Channel* into low-level commands. We do not deem our requirements for the robot platforms too high: the robot should provide some controllers to move either the end-effector(s) and the base in Cartesian space, the joints, or combinations of both.

The action command structure in TeleMoMa relayed to the *Robot Interface* can either contain values in task-space (end-effector Cartesian relative motion), joint space (e.g., torso commands or motion to other joints) and/or velocities (e.g., base commands), or different combinations of those, as specified by the user during deployment. The *Robot Interface* processes these commands based on the particular robot embodiment, filters out the unusable action components (such as left hand commands for a single-armed robot like Fetch), and maps the rest to the robot using the preferred choices of controllers such as operational space control [1] to control one task frame, or whole-body control [2, 33] to command the entire robot jointly, or having separate controls for each part of the robot.

IV. EXPERIMENTS

In our experiments we seek to answer the following questions: (1) What are the benefits of TeleMoMa’s modularity? (Sec. IV-A) (2) Can TeleMoMa collect high-quality data for imitation learning? (Sec. IV-B), (3) How does TeleMoMa perform in remote teleoperation of the robot with possible network delays? (Sec. IV-C), and (4) What is the effect of different robot embodiments and the gap between simulation and real in the usability of TeleMoMa? (Sec. IV-D).

A. User Study

To assess the performance of different teleoperation modalities in the TeleMoMa framework, we performed two user studies with the PAL Tiago++ robot. We compared three teleoperation modalities described in Sec. III-A: *VR*, in which the user controls the robot’s arms with the Oculus controllers and the base and torso with the controller joysticks; *Vision*, in which the user’s pose is tracked with an RGB-D camera to control the arms, torso and base motion; and *VR+Vision* combining both modalities, in which the robot’s arms are controlled using the Oculus controllers and the base and torso motion is controlled via human pose tracking from RGB-D data.

In the first user study, we compared the three modalities (*VR*, *Vision*, *VR + Vision*) to assess the completion time in two tasks: *cover table* (Fig. 3(a)), in which the robot must grasp a tablecloth with both hands and drape it over a table, and *dusting* (Fig. 3(j)), in which the robot must dust a table with books resting on top. Both tasks, but especially the *dusting* task, benefit from the simultaneous motion of base and arm(s), i.e., whole-body motion, as enabled by TeleMoMa since the robot is required to navigate around the desk while periodically moving the hands to clear out any dust.

We recruited 12 participants with varying levels of teleoperation experience. Each user was given the same instructions and a brief practice period with each modality. The

order in which users received the devices was randomized. The completion times for successful trials are provided in Fig. 4. The only failures observed occurred with the *Vision* modality (3 fails out of 12 *dusting* trials) due to noise and inaccuracies in the pose tracking. We observe that in the *cover table* task, performance is comparable across teleoperation modalities. However, in the *dusting* task, pure *VR* is generally slower than *VR + Vision* or *Vision* alone due to the lack of intuitive whole-body teleoperation: because moving the base requires using the joysticks on the controllers, users tended to only move the arm or the base one at a given time. The results indicate that on their own, both *VR* and *Vision* present drawbacks pertaining to their individual modalities, but when combined in the form of *VR + Vision*, TeleMoMa can overcome their individual drawbacks to enable an improved teleoperation experience. These results support empirically the importance of enabling multiple input modalities and their combination for teleoperation of mobile manipulators, and TeleMoMa’s potential for enabling data collection in more complex mobile manipulation tasks beyond pick and place.

In the second study, we sought to assess whether TeleMoMa users improve over time by measuring their learning curve. We recruited 6 participants with varying levels of teleoperation experience and compared two modalities (*VR* and *VR + Vision*, order randomized) on the *pick up* task (Fig. 3(i)). In this task, the robot must lower its torso in order to grasp a towel from the floor, hand the towel from one hand to the other, navigate to a table, and place the towel on the table. Users completed three consecutive trials with each modality and completion times were recorded. The results are visualized in Fig. 5. We observe that new users generally improve at completing tasks with the system, with an average decrease of 29% and 26% in task completion time over three trials for the *VR* and *VR + Vision* respectively. The completion times were generally similar between *VR* and *VR + Vision*, with some slower times in the *VR + Vision* modality owing to the increased difficulty of controlling additional degrees of freedom simultaneously and the additional noise introduced by the vision-based human interface. Despite slowing performance in some trials of this task, the additional capabilities from *VR + Vision* have the potential to unlock new whole-body control applications once mastered. Taken together, these two user studies demonstrate the benefits of TeleMoMa as a modular teleoperation system.

B. Imitation Learning with TeleMoMa’s Data

To empirically evaluate the quality of the data collected with TeleMoMa, we train several visuomotor policies with behavioral cloning [44] using the data collected on a Tiago++ robot (*real*). We consider three diverse mobile manipulation tasks:

- *cover table*: Similar to the one described in Sec. IV-A, the tasks involves bimanual grasping of a tablecloth and draping it over a table (Fig. 3(a)).
- *slide chair*: A bimanual task, that requires the robot to navigate and align itself behind a chair, grasp it, and

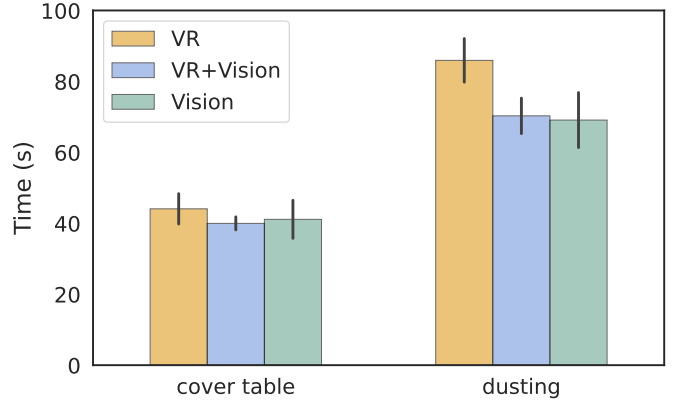


Fig. 4: **User Study 1: Completion Time.** Vision modalities outperform only-VR for the more challenging *dusting* task. Error bars denote the standard error of the mean.

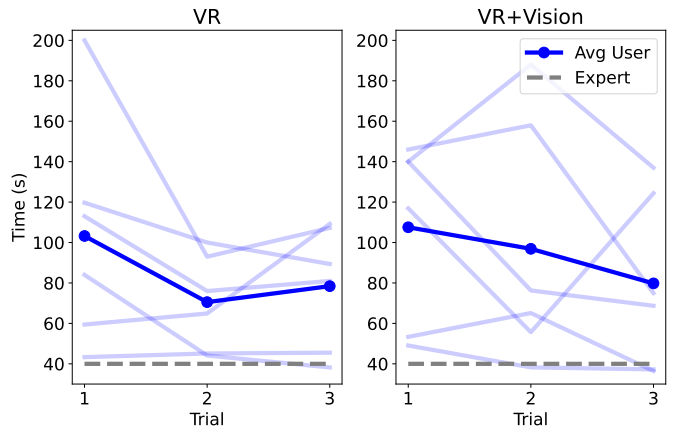


Fig. 5: **User Study 2: User Improvement, Learning Curve.** New users generally improve at completing the *pick up* task with TeleMoMa across teleoperation modalities. Transparent lines show individual learning curves.

push the chair under a table (Fig. 3(c)).

- *serve bread*: In a real kitchen setting, the robot is required to navigate to the kitchen counter, pick a bag of bread, and deliver it to the breakfast table (Fig. 3(e)).

We collected 50 demonstrations each for *slide chair* and *serve bread* tasks and 100 demonstrations for *cover table* task using the combined *VR + Vision* interface of TeleMoMa. Additional demonstrations in the *cover table* were necessary to allow the policies to learn the necessary accurate grasps on the cloth.

Policy Architecture, Observations and Actions. We used a feed-forward MLP (BC) and a recurrent LSTM based network (BC-RNN) [32] with a sequence length of 10. The inputs to all policies included RGB-D images obtained from two realsense cameras attached on each shoulder of the robot, end-effector poses of the hands, gripper state, and the change in the mobile base pose obtained from the odometry of the robot. The policies output a 17-dimensional action space: 6D Cartesian deltas and a gripper command for each of the hands, and linear and angular velocities for the base.

Comparing Input Modalities. To analyze the importance of depth sensing in learning mobile manipulation tasks, we

TABLE II: Performance between IL policies trained with RGB vs. RGBD images as inputs. Successes measured over 10 rollouts.

Modality	Cover Table		Slide Chair		Serve Bread	
	RGB	RGB-D	RGB	RGB-D	RGB	RGB-D
BC	60	60	40	60	20	40
BC-RNN	70	90	50	80	30	70

TABLE III: IL Policy performance scale with data. Successes measured over 10 rollouts.

Fraction of data	Cover Table		Slide Chair		Serve Bread	
	50%	100%	50%	100%	50%	100%
BC	60	60	40	60	30	40
BC-RNN	60	90	70	80	40	70

train two sets of policies: the first set was trained exclusively on RGB observations, while the second combined RGB and Depth. The performance of the two sets of policies for each of the tasks is summarized in Table II. Our analysis reveals a consistent trend: irrespective of the policy architecture, the inclusion of depth information markedly enhances performance across all tasks. Qualitatively, we observe that policies trained using depth can position the base better, significantly improving the efficacy of subsequent arm actions. These findings suggest that depth information is a crucial component for the development of effective mobile manipulation policies, and that the strong dependency between base and arm actions is one of the main challenges in IL for mobile manipulation.

Performance with Different Amounts of Data. To investigate how data volume influences policy performance, we experimented with two distinct policy groups: the first group was trained using the complete dataset we gathered for each task, while the second group utilized only 50% of these collected demonstrations. The results are summarized in Table III; we observe that policies trained with the full dataset consistently outperform those trained with half the data, demonstrating the importance of dataset size in imitation learning, especially in this low-data regime. We additionally notice that BC-RNN strictly outperforms regular BC in all tasks, demonstrating the significance of temporal dependencies for learning mobile manipulation tasks.

In general, the above experiments provide compelling evidence that IL policies trained with data collected using TeleMoMa can reliably perform complex mobile manipulation tasks, thus indicating that TeleMoMa can facilitate high-quality data collection for imitation learning. We demonstrate more imitation results in the sim environment in Appendix C.

C. Remote Teleoperation

TeleMoMa’s architecture allows a remote demonstrator to control the robot from a client computer connected over the internet. Instead of watching the robot on-site, the demonstrator is provided with camera streams transmitted by the teleoperation channel from the robot’s onboard sensors. To

minimize communication delays, TeleMoMa 1) sends compressed sensor images from the robot and decompresses them on the client, and 2) in the case of a vision-based human interface, TeleMoMa processes the RGB-D images from the vision interface on the client side and only sends the action commands over the teleoperation channel. For other interfaces, the demonstrated action commands are directly sent to the TeleMoMa’s robot interface.

We demonstrate the remote teleoperation capability of TeleMoMa on several combinations of robot hardware and user interfaces. To evaluate the effects of communication delays, we compare the task completion time between on-site and remote demonstrations using Tiago++ and Toyota HSR each on two different tasks. The `cover table` and the `slide chair` tasks are completed using Tiago++ with the on-site `VR + Vision` interface and three remote interfaces (`VR`, `Vision`, `VR + Vision`). The `re-shelve chips` task, in which the robot must move the misplaced chips to the lower shelf (Fig. 3(h)), and the `open fridge` task, in which the robot must open a fridge (Fig. 3(g)), are completed using HSR with the `Vision` interface. The demonstrations are provided by an expert user of each robot. The Wi-Fi speed is about 100 Mbps as measured on the HSR. Fig. 6(a) shows the completion time in each modality averaged over 3 runs. We observe that remote human demonstrators have slower reaction times due to delays and limited resolutions of the camera streams, but TeleMoMa provides the capability to successfully complete the tasks under regular network conditions.

D. Comparing Different Embodiments and Sim vs. Real

In the final set of experiments, we seek to study how the domain (sim vs. real) and the type of robot (Tiago vs. HSR and Tiago-sim vs. Fetch-sim) influence the teleoperation behavior for the same tasks.

1) *Sim vs. Real:* Fig. 6(b) depicts the results of comparing completion time for `cover table` and `slide chair` tasks in simulation and real environment using a Tiago robot. We use `sim` time for simulation evaluation because of OmniGibson’s sub-realtime soft-body simulation. By maintaining consistency across the robot, the task, and the teleoperation interface, we find that for both tasks the completion time in simulation and real are close, demonstrating that the simulation environment in OmniGibson is a good proxy for mobile manipulation in the real world, and that teleoperating with TeleMoMa provides a natural mechanism to collect demonstrations in sim.

2) *Comparing Embodiments:* We additionally compare how the completion time varies as we change the robot being teleoperated by maintaining the task, teleoperation interface and reality to be consistent. We compare Tiago and HSR on `re-shelve chips` and `open fridge` tasks and depict the results in Fig. 6(c, right). We observe that the higher number of degrees of freedom offered by Tiago compared to HSR allows more fluid motion during teleoperation and enables a more efficient (faster) completion of the task.

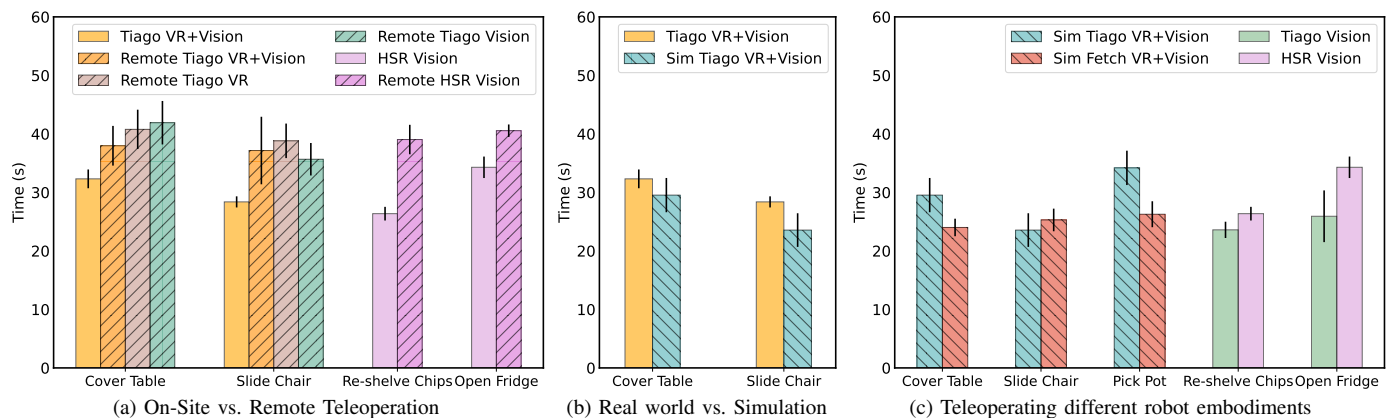


Fig. 6: Completion times in different experiments with TeleMoMa. The bars indicate the mean and standard deviation of several trials (see text). *From left to right:* Comparing completion times for tasks performed on-site and remote, with HSR and Tiago; Completion times for real vs. simulated tasks with Tiago; Completion times for different robot embodiments on the same tasks in the real world and simulation. TeleMoMa allows for multiple tasks in simulation and the real world, with several embodiments

In simulation, we compare Tiago and Fetch on `cover table`, `slide chair`, and `pick pot` tasks and depict the results in Fig. 6(c, left). For the `pick pot` task, we enabled sticky grasping (creating a controllable constraint between hand and object) since the task would be infeasible otherwise for a single-armed robot like Fetch. We observe that Fetch is faster than Tiago on tasks requiring table-top manipulations, possibly due to Fetch’s larger size and longer arms, making manipulation easier for users.

V. CONCLUSIONS

We presented TeleMoMa, a novel teleoperation system for mobile manipulators that enables versatility through modularity. While no single teleoperation interface provides all benefits of enabling dexterous, whole-body mobile teleoperation while remaining low cost and scalable, our general, modular teleoperation interface provides the ability to combine multiple existing modalities combining also some of the benefits of them. This results in a performant data collection system that scales to many different robots and tasks, as indicated by our user studies, imitation learning, remote teleoperation, and comparisons between embodiments and sim vs. real. We note some limitations of TeleMoMa. First, when tracking human pose from RGB data, noise and inaccuracies can impact a user’s ability to accomplish tasks. Combining vision with a more accurate interface like VR enables accurate arm control and synchronization of base and arm movement, but would still benefit from better visual pose-tracking models. Second, occlusion presents a challenge for the vision-based modalities, as the camera placement has an impact on the operator’s visibility of the robot’s workspace. This can be mitigated by carefully choosing a camera placement, using multiple cameras, or rendering robot observations on a screen. Extending TeleMoMa to incorporate a puppeteering human interface would enable even more accurate tasks at the cost of mobility.

In closing, we have demonstrated TeleMoMa, a general, modular, accessible teleoperation system that enables collec-

tion of high-quality expert demonstration data for a variety of complex and novel mobile manipulation tasks. We showed TeleMoMa’s generality by teleoperating multiple different robots in simulation and reality, and conducted user studies to verify the usability of the system’s various modalities. We hope that our system lowers the barrier of entry for researchers to collect high-quality demonstrations for mobile manipulation, and helps unlock new mobile manipulation capabilities.

REFERENCES

- [1] A unified approach for motion and force control of robot manipulators: The operational space formulation. *IEEE Journal on Robotics and Automation*, 3(1):43–53, 1987.
- [2] Whole-body dynamic behavior and control of human-like robots. *International Journal of Humanoid Robotics*, 1(01):29–43, 2004.
- [3] Miguel Arduengo, Ana Arduengo, Adrià Colomé, Joan Lobo-Prat, and Carme Torras. Human to robot whole-body motion transfer. In *2020 IEEE-RAS 20th International Conference on Humanoid Robots (Humanoids)*, pages 299–305. IEEE, 2021.
- [4] Brenna D Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. A survey of robot learning from demonstration. *Robotics and autonomous systems*, 57(5):469–483, 2009.
- [5] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- [6] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.
- [7] Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho, Julian Ibarz,

- Alex Irpan, Eric Jang, Ryan Julian, et al. Do as i can, not as i say: Grounding language in robotic affordances. In *Conference on Robot Learning*, pages 287–318. PMLR, 2023.
- [8] Yevgen Chebotar, Quan Vuong, Karol Hausman, Fei Xia, Yao Lu, Alex Irpan, Aviral Kumar, Tianhe Yu, Alexander Herzog, Karl Pertsch, et al. Q-transformer: Scalable offline reinforcement learning via autoregressive q-functions. In *Conference on Robot Learning*, pages 3909–3928. PMLR, 2023.
- [9] Open X-Embodiment Collaboration. Open X-Embodiment: Robotic learning datasets and RT-X models. <https://arxiv.org/abs/2310.08864>, 2023.
- [10] Shivin Dass, Karl Pertsch, Hejia Zhang, Youngwoon Lee, Joseph J Lim, and Stefanos Nikolaidis. Pato: Policy assisted teleoperation for scalable robot data collection. *arXiv preprint arXiv:2212.04708*, 2022.
- [11] Joseph DelPreto, Jeffrey I Lipton, Lindsay Sanneman, Aidan J Fay, Christopher Fourie, Changhyun Choi, and Daniela Rus. Helping robots learn: a human-robot master-apprentice model using demonstrations via virtual reality teleoperation. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 10226–10233. IEEE, 2020.
- [12] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palme: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.
- [13] 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.
- [14] Roya Firoozi, Johnathan Tucker, Stephen Tian, Anirudha Majumdar, Jiankai Sun, Weiyu Liu, Yuke Zhu, Shuran Song, Ashish Kapoor, Karol Hausman, et al. Foundation models in robotics: Applications, challenges, and the future. *arXiv preprint arXiv:2312.07843*, 2023.
- [15] Lars Fritsche, Felix Unverzag, Jan Peters, and Roberto Calandra. First-person tele-operation of a humanoid robot. In *2015 IEEE-RAS 15th International Conference on Humanoid Robots (Humanoids)*, pages 997–1002. IEEE, 2015.
- [16] 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.
- [17] Bryan R Galarza, Paulina Ayala, Santiago Manzano, and Marcelo V Garcia. Virtual reality teleoperation system for mobile robot manipulation. *Robotics*, 12(6):163, 2023.
- [18] C Gan, J Schwartz, S Alter, M Schrimpf, J Traer, J De Freitas, J Kubilius, A Bhandwaldar, N Haber, M Sano, et al. Threedworld: A platform for interactive multi-modal physical simulation. *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [19] Xiaofeng Gao, Ran Gong, Tianmin Shu, Xu Xie, Shu Wang, and Song-Chun Zhu. Vrkitcchen: an interactive 3d virtual environment for task-oriented learning. *arXiv preprint arXiv:1903.05757*, 2019.
- [20] Alberto Garcia-Garcia, Pablo Martinez-Gonzalez, Sergiu Oprea, John Alejandro Castro-Vargas, Sergio Orts-Escolano, Jose Garcia-Rodriguez, and Alvaro Jover-Alvarez. The robotrix: An extremely photorealistic and very-large-scale indoor dataset of sequences with robot trajectories and interactions. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 6790–6797. IEEE, 2018.
- [21] Ankur Handa, Karl Van Wyk, Wei Yang, Jacky Liang, Yu-Wei Chao, Qian Wan, Stan Birchfield, Nathan Ratliff, and Dieter Fox. Dexpivot: Vision-based teleoperation of dexterous robotic hand-arm system. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9164–9170. IEEE, 2020.
- [22] Ryan Hoque, Lawrence Yunliang Chen, Satvik Sharma, Karthik Dharmarajan, Brijen Thananjeyan, Pieter Abbeel, and Ken Goldberg. Fleet-dagger: Interactive robot fleet learning with scalable human supervision. In *Conference on Robot Learning*, pages 368–380. PMLR, 2023.
- [23] Gayane Kazhoyan, Alina Hawkin, Sebastian Koralewski, Andrei Haidu, and Michael Beetz. Learning motion parameterizations of mobile pick and place actions from observing humans in virtual environments. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 9736–9743. IEEE, 2020.
- [24] Jens Kober, J. Andrew Bagnell, and Jan Peters. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11):1238–1274, September 2013. ISSN 0278-3649, 1741-3176. doi: 10.1177/0278364913495721. URL <http://journals.sagepub.com/doi/10.1177/0278364913495721>.
- [25] Franziska Krebs, Andre Meixner, Isabel Patzer, and Tamim Asfour. The kit bimanual manipulation dataset. In *2020 IEEE-RAS 20th International Conference on Humanoid Robots (Humanoids)*, pages 499–506. IEEE, 2021.
- [26] Christian Lenz and Sven Behnke. Bimanual telemanipulation with force and haptic feedback through an anthropomorphic avatar system. *Robotics and Autonomous Systems*, 161:104338, 2023.
- [27] Chengshu Li, Ruohan Zhang, Josiah Wong, Cem Gokmen, Sanjana Srivastava, Roberto Martín-Martín, Chen Wang, Gabriel Levine, Michael Lingelbach, Jiankai Sun, et al. Behavior-1k: A benchmark for embodied ai with 1,000 everyday activities and realistic simulation. In *Conference on Robot Learning*, pages 80–93. PMLR, 2023.
- [28] Jeffrey I Lipton, Aidan J Fay, and Daniela Rus. Baxter’s homunculus: Virtual reality spaces for teleoperation in manufacturing. *IEEE Robotics and Automation Letters*,

- 3(1):179–186, 2017.
- [29] Camillo Lugaresi, Jiuqiang Tang, Hadon Nash, Chris McClanahan, Esha Uboweja, Michael Hays, Fan Zhang, Chuo-Ling Chang, Ming Yong, Juhyun Lee, et al. Medi-pipe: A framework for perceiving and processing reality. In *Third workshop on computer vision for AR/VR at IEEE computer vision and pattern recognition (CVPR)*, volume 2019, 2019.
- [30] 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.
- [31] Ajay Mandlekar, Jonathan Booher, Max Spero, Albert Tung, Anchit Gupta, Yuke Zhu, Animesh Garg, Silvio Savarese, and Li Fei-Fei. Scaling robot supervision to hundreds of hours with roboturk: Robotic manipulation dataset through human reasoning and dexterity. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1048–1055. IEEE, 2019.
- [32] Ajay Mandlekar, Danfei Xu, Josiah Wong, Soroush Nasiriany, Chen Wang, Rohun Kulkarni, Li Fei-Fei, Silvio Savarese, Yuke Zhu, and Roberto Martín-Martín. What matters in learning from offline human demonstrations for robot manipulation. In *arXiv preprint arXiv:2108.03298*, 2021.
- [33] Nicolas Mansard, Olivier Stasse, Paul Evrard, and Abderrahmane Kheddar. A versatile generalized inverted kinematics implementation for collaborative working humanoid robots: The stack of tasks. In *2009 International conference on advanced robotics*, pages 1–6. IEEE, 2009.
- [34] Pablo Martinez-Gonzalez, Sergiu Oprea, Alberto Garcia-Garcia, Alvaro Jover-Alvarez, Sergio Orts-Escolano, and Jose Garcia-Rodriguez. Unrealrox: an extremely photorealistic virtual reality environment for robotics simulations and synthetic data generation. *Virtual Reality*, 24:271–288, 2020.
- [35] Yutaro Matsuura, Kento Kawaharazuka, Naoki Hiraoka, Kunio Kojima, Kei Okada, and Masayuki Inaba. Development of a whole-body work imitation learning system by a biped and bi-armed humanoid. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 10374–10381. IEEE, 2023.
- [36] Jordi Pages, Luca Marchionni, and Francesco Ferro. Tiago: the modular robot that adapts to different research needs. In *International workshop on robot modularity, IROS*, volume 290, 2016.
- [37] Luigi Penco, Kazuhiko Momose, Stephen McCrory, Dexter Anderson, Nicholas Kitchel, Duncan Calvert, and Robert J Griffin. Mixed reality teleoperation assistance for direct control of humanoids. *IEEE Robotics and Automation Letters*, 2024.
- [38] Amartya Purushottam, Christopher Xu, Yeongtae Jung, and Joao Ramos. Dynamic mobile manipulation via whole-body bilateral teleoperation of a wheeled humanoid. *IEEE Robotics and Automation Letters*, 2023.
- [39] Yuzhe Qin, Wei Yang, Binghao Huang, Karl Van Wyk, Hao Su, Xiaolong Wang, Yu-Wei Chao, and Dieter Fox. Anyteleop: A general vision-based dexterous robot arm-hand teleoperation system. *arXiv preprint arXiv:2307.04577*, 2023.
- [40] Rouhollah Rahmatizadeh, Pooya Abolghasemi, Ladislau Bölöni, and Sergey Levine. Vision-based multi-task manipulation for inexpensive robots using end-to-end learning from demonstration. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 3758–3765. IEEE, 2018.
- [41] Ellis Ratner, Benjamin Cohen, Mike Phillips, and Maxim Likhachev. A web-based infrastructure for recording user demonstrations of mobile manipulation tasks. In *2015 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5523–5530. IEEE, 2015.
- [42] Harish Ravichandar, Athanasios S Polydoros, Sonia Chernova, and Aude Billard. Recent advances in robot learning from demonstration. *Annual review of control, robotics, and autonomous systems*, 3:297–330, 2020.
- [43] Dongseok Ryu, Jae-Bok Song, Changhyun Cho, Sungchul Kang, and Munsang Kim. Development of a six dof haptic master for teleoperation of a mobile manipulator. *Mechatronics*, 20(2):181–191, 2010.
- [44] Stefan Schaal. Is imitation learning the route to humanoid robots? *Trends in cognitive sciences*, 3(6):233–242, 1999.
- [45] Max Schwarz, Christian Lenz, Raphael Memmesheimer, Bastian Pätzold, Andre Rochow, Michael Schreiber, and Sven Behnke. Robust immersive telepresence and mobile telemanipulation: Nimbro wins ana avatar xprize finals. *arXiv preprint arXiv:2303.03297*, 2023.
- [46] Mingyo Seo, Steve Han, Kyutae Sim, Seung Hyeon Bang, Carlos Gonzalez, Luis Sentis, and Yuke Zhu. Deep imitation learning for humanoid loco-manipulation through human teleoperation. In *2023 IEEE-RAS 22nd International Conference on Humanoid Robots (Humanoids)*, pages 1–8. IEEE, 2023.
- [47] Adam Setapen, Michael Quinlan, and Peter Stone. Marionet: Motion acquisition for robots through iterative online evaluative training. In *Ninth International Conference on Autonomous Agents and Multiagent Systems - Agents Learning Interactively from Human Teachers Workshop (AAMAS - ALIHT)*, May 2010.
- [48] Bruno Siciliano, Oussama Khatib, and Torsten Kröger. *Springer handbook of robotics*, volume 200. Springer, 2008.
- [49] Aravind Sivakumar, Kenneth Shaw, and Deepak Pathak. Robotic telekinesis: Learning a robotic hand imitator by watching humans on youtube. *arXiv preprint arXiv:2202.10448*, 2022.
- [50] Christopher Stanton, Anton Bogdanovych, and Edward Ratanasena. Teleoperation of a humanoid robot using full-body motion capture, example movements, and machine learning. In *Proc. Australasian Conference on Robotics and Automation*, volume 8, page 51, 2012.

- [51] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Charles Xu, Jianlan Luo, et al. Octo: An open-source generalist robot policy, 2023.
- [52] Albert Tung, Josiah Wong, Ajay Mandlekar, Roberto Martín-Martín, Yuke Zhu, Li Fei-Fei, and Silvio Savarese. Learning multi-arm manipulation through collaborative teleoperation. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9212–9219. IEEE, 2021.
- [53] Jean Vertut and Philippe Coiffet. *Teleoperations and robotics: evolution and development*. Prentice-Hall, Inc., 1986.
- [54] Homer Walke, Kevin Black, Abraham Lee, Moo Jin Kim, Max Du, Chongyi Zheng, Tony Zhao, Philippe Hansen-Estruch, Quan Vuong, Andre He, Vivek Myers, Kuan Fang, Chelsea Finn, and Sergey Levine. Bridgedata v2: A dataset for robot learning at scale. In *Conference on Robot Learning (CoRL)*, 2023.
- [55] Baocheng Wang, Zhijun Li, Wenjun Ye, and Qing Xie. Development of human-machine interface for teleoperation of a mobile manipulator. *International Journal of Control, Automation and Systems*, 10:1225–1231, 2012.
- [56] David Whitney, Eric Rosen, Daniel Ullman, Elizabeth Phillips, and Stefanie Tellex. Ros reality: A virtual reality framework using consumer-grade hardware for ros-enabled robots. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1–9. IEEE, 2018.
- [57] Melonee Wise, Michael Ferguson, Derek King, Eric Diehr, and David Dymesich. Fetch and freight: Standard platforms for service robot applications. In *Workshop on autonomous mobile service robots*, pages 1–6, 2016.
- [58] Josiah Wong, Albert Tung, Andrey Kurenkov, Ajay Mandlekar, Li Fei-Fei, Silvio Savarese, and Roberto Martín-Martín. Error-aware imitation learning from teleoperation data for mobile manipulation. In *Conference on Robot Learning*, pages 1367–1378. PMLR, 2022.
- [59] 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.
- [60] Takashi Yamamoto, Koji Terada, Akiyoshi Ochiai, Fuminori Saito, Yoshiaki Asahara, and Kazuto Murase. Development of human support robot as the research platform of a domestic mobile manipulator. *ROBOMECH journal*, 6(1):1–15, 2019.
- [61] Taozheng Yang, Ya Jing, Hongtao Wu, Jiafeng Xu, Kuankuan Sima, Guangzeng Chen, Qie Sima, and Tao Kong. Moma-force: Visual-force imitation for real-world mobile manipulation. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 6847–6852. IEEE, 2023.
- [62] 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 *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5628–5635. IEEE, 2018.
- [63] Zhijun Zhang, Yaru Niu, Ziyi Yan, and Shuyang Lin. Real-time whole-body imitation by humanoid robots and task-oriented teleoperation using an analytical mapping method and quantitative evaluation. *Applied Sciences*, 8(10):2005, 2018.
- [64] 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.

APPENDIX

A. Criterion for Table 1

We provide a detailed explanation for each column of Table 1, including the criteria used to categorize methods.

- Teleoperation Support

- 1) *Cost / Accessibility*: We identified three tiers of price based on commercially available systems or disclosed cost.

\$: \$0 – 1,000 (VR, Vision, Phone)
 \$\$: \$1,000 – 10,000 (Mocap Systems)
 \$\$\$: \$10,000+ (Custom Hardware)

- 2) *Modular*: True if the method is modular in the sense that it supports multiple input modalities or combinations thereof. TeleMoMa is the only method that meets this criteria.
- 3) *Modality*: Modality describes the human interface used for teleoperation (e.g. virtual reality (VR), puppeteering with a kinematically similar device, motion capture systems (Mocap), etc.).

- Robot Support

- 1) *Bimanual*: True if the paper demonstrates bimanual teleoperation.
- 2) *Height Control*: True if the paper demonstrates control of the robot’s torso joint.
- 3) *Whole-Body Teleoperation*: True if simultaneous arm and base motion is enabled by the method.
- 4) *Robot Agnostic*: True if the method works for many different robots; false if it is specific to a particular platform.
- 5) *Action Space*: “EE Pose(s)” denotes control of the robot’s end-effector(s) in Cartesian space, whereas “Joint Pos.” indicates joint-space control for the arms and/or torso. Base Vel. indicates control of the base velocity; TRILL [46] allows users to select among predefined gaits with a VR controller, denoted “Gait”. MOMA-Force enables teleoperation of end-effector Cartesian pose through kinesthetic teaching and additionally records desired end-effector wrenches, denoted “EE Pose and Wrench”. TeleMoMa allows users to control end-effector Cartesian pose, base velocity, and torso joint position; it is also readily extensible to joint control when tracking human pose, but this is left for future work.

B. Method Details

Following we describe how TeleMoMa facilitates the use of mobile phones, spacemouse, and keyboards as part of its *Human Interface* (Sec. IV-A). We are also open-sourcing the code to the community to facilitate plug-and-play teleoperation for mobile manipulators to improve data collection efficiency.

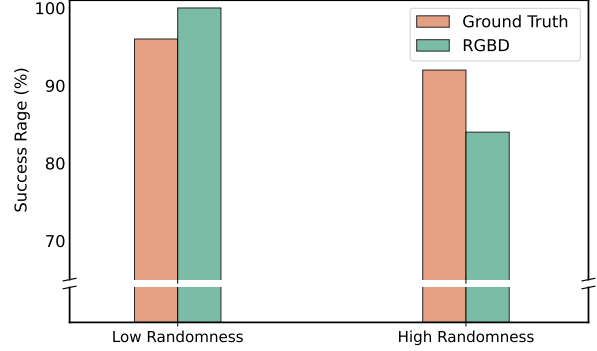


Fig. 7: **IL Results in Simulation.** Policy with RGBD input yields comparable performance to policy with ground truth chair positions as input.

- 1) *Mobile Phone*: We created an app using the ARKit development kit to track the position and orientation of the mobile phone, which sends commands over the network. Similar to *Virtual Reality Controller* (Sec. IV-A2), the end effector is commanded in the task space and the relative pose displacement per frame of the mobile phone is calculated and mapped to the robot end effector. The gripper is controlled by dedicated buttons in the mobile app. Additionally, simultaneous control of left and right arms can be facilitated if two mobile phones are running the app, each phone controlling one of the arms. Mobile phones currently don’t support navigation capabilities, but can be combined with other modalities such as the *Vision-based Human Interface* (Sec. IV-A1) to facilitate mobile base movements.

- 2) *Spacemouse*: Spacemouse has only 6-degrees of freedom, which is why we use mode switching, and control each part of the robot independently. The users can switch modes by pressing one of the side buttons of the spacemouse and switch between controlling left arm, right arm, base and torso. Two spacemouse’ can also be used simultaneously for controlling each of the arms and minimizing the mode switching. The displacement of the spacemouse in each of the 6 degrees of freedom is tracked and sent as the delta commands to control the arms. For the base and torso, only the required displacements are used to send commands, while the remaining ones are discarded. The gripper can be toggled by pressing the remaining side button when the spacemouse mode is controlling the corresponding arm. Spacemouse gains significantly from modularity offered by TeleMoMa, by minimizing mode switching thus gaining more fluid control of the robot.

- 3) *Keyboard*: Keyboard presses are asynchronously read by the device listeners and each key is mapped to a single DoF of the mobile manipulator. Each key increases / decreases one of the DoFs in the Cartesian space by some preset amount. This results in a large number of keys that the teleoperator has to remember for controlling the robot. Instead, using a smaller set of keys for controlling for instance, just the base, while controlling arms with something more intuitive such as the

Hyperparameters	Value
Behavior Cloning (BC)	
train steps (x500)	500
batch size	32
optimizer	Adam
learning rate	1e-4
image & depth encoder	resnet-18
policy (w x d)	512x2
action parameterization	GMM
Recurrent BC (BC-RNN)	
train steps (x500)	500
batch size	16
optimizer	Adam
learning rate	1e-4
image & depth encoder	resnet-18
LSTM hidden dim	1000
LSTM num. layers	2
skill horizon	10
action parameterization	GMM

TABLE IV: Hyperparameters for the imitation policies (the hyperparameter values were kept consistent across tasks)

spacemouse can drastically improve the teleoperation experience on both the interfaces, minimizing the mode switching in case of spacemouse, and reducing the number of keys to keep track of on the keyboard.

C. Imitation Results in Simulation

We show the imitation results of the `slide chair` task in simulation here. We collected 100 demos in OmniGibson, and trained 2 policies using BC with different input observations: one with RGB-D image from the head camera, and the other with oracle chair positions in both world frame and robot base frame from the simulation environment. Robot proprioception, including end effector poses for two arms in base frame, and the base position and velocity in world frame, are also provided as observation input. We evaluated the policy on two task configurations: first with low randomness, where the chair position is uniformly sampled within 0.2 meters parallel to the robot, and second with high randomness, where the sampling interval is 1 meters. Each policy is evaluated with 25 rollouts under these conditions.

The results are shown in Fig. 7. We observed that, the performance of policies under high randomness is worse than under low randomness, which is expected because of the increased difficulty. We additionally observe that in both low and high randomness settings, policy trained with RGB-D input performs comparable to the one trained with ground truth chair positions, indicating that the policies are able to extract meaningful environment specific details from images and depth. Qualitatively, we observe that the causes of failure include misalignment between the robot and the chair, slippage of robot grippers, and knocking over the chair due to the application of excessive force.

D. Imitation Learning Policy Hyperparameters

We performed imitation learning on one simulated (`slide chair` – Appendix Sec. C) and three real world tasks – `cover table` (Fig. 3(a)), `slide chair` (Fig. 3(b)) and `serve bread` (Fig. 3(c)), that require synchronized hand

and base motions. The results and their analysis are presented in Sec. V-B. We used RoboMimic [32] for training the policies. Comprehensive details of the policy architecture and hyperparameters used for training are provided in Table IV. Note that the same hyperparameters were used across all tasks, and across simulation and real environments.

Furthermore, in an effort to facilitate and encourage ongoing research in mobile manipulation, the dataset collected on all the tasks will be made available along with the code.