# Humanoid Manipulation Interface: Humanoid Whole-Body Manipulation from Robot-Free Demonstrations

Ruiqian Nai*[1,2], Boyuan Zheng*[1,2], Junming Zhao*[1,2], Haodong Zhu[1], Sicong Dai[1,†],
Zunhao Chen[1], Yihang Hu[1,2], Yingdong Hu[1,2], Tong Zhang[1,2], Chuan Wen[3], Yang Gao[1,2]
*Equal contribution, [1]Tsinghua University, [2]Shanghai Qi Zhi Institute,
[3]Shanghai Jiao Tong University
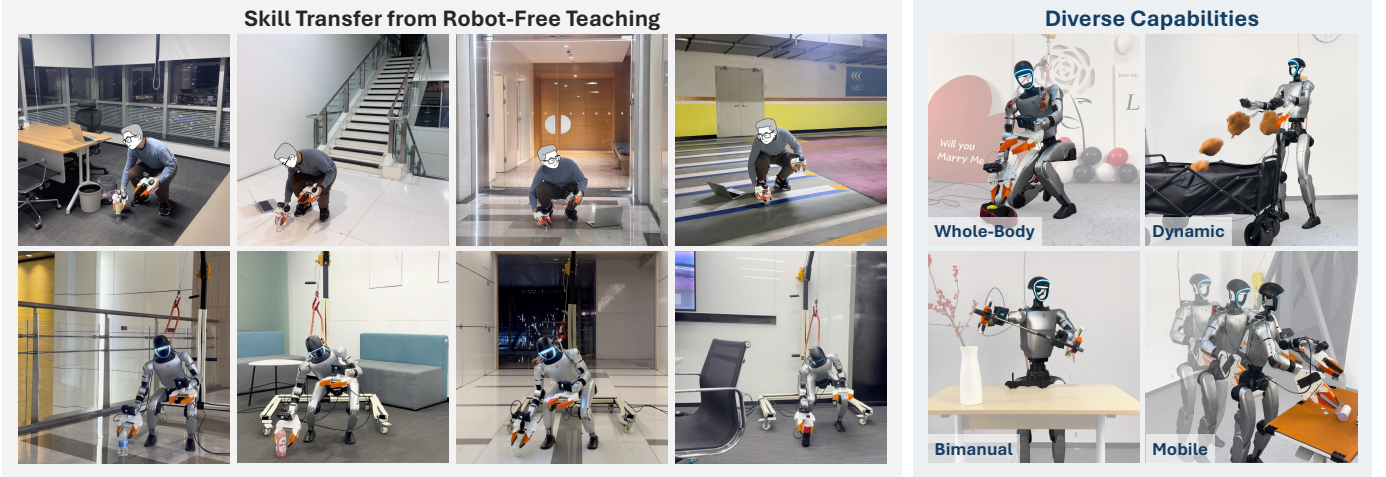https://humanoid-manipulation-interface.github.io

Fig. 1: **Humanoid Manipulation Interface (HuMI).** **Left:** Our portable, robot-free data collection facilitates skill transfer from human to humanoid across diverse, unstructured environments. **Right:** The framework enables a wide repertoire of complex whole-body behaviors.

*Abstract*—Current approaches for humanoid whole-body manipulation, primarily relying on teleoperation or visual sim-to-real reinforcement learning, are hindered by hardware logistics and complex reward engineering. Consequently, demonstrated autonomous skills remain limited and are typically restricted to controlled environments. In this paper, we present the Humanoid Manipulation Interface (HuMI), a portable and efficient framework for learning diverse whole-body manipulation tasks across various environments. HuMI enables robot-free data collection by capturing rich whole-body motion using portable hardware. This data drives a hierarchical learning pipeline that translates human motions into dexterous and feasible humanoid skills. Extensive experiments across five whole-body tasks—including kneeling, squatting, tossing, walking, and bimanual manipulation—demonstrate that HuMI achieves a 3x increase in data collection efficiency compared to teleoperation and attains a 70% success rate in unseen environments.

## I. INTRODUCTION

Humans expertly coordinate their entire bodies for manipulation, whether squatting to retrieve objects or bending to reach low tables. With their high degrees of freedom, humanoid robots are expected to exhibit similar whole-body capabilities,

tightly coordinating all joints for manipulation using onboard perception.

To achieve this, recent research employs visual sim-to-real reinforcement learning (RL) [15, 25, 46] or imitation learning from teleoperation [2, 21, 22, 52, 54]. However, these methods are labor-intensive: RL demands diverse assets and meticulous reward engineering, while teleoperation requires significant expertise to manage balance and controller inaccuracies. Consequently, current methods demonstrate few autonomous tasks in fixed lab environments [2, 15, 21, 25, 46, 54]. These tasks exhibit limited whole-body coordination, typically restricting robots to upright walking combined with simple actions like transporting objects, opening doors, or kicking boxes.

In this project, our goal is to enable humanoid robots to perform diverse tasks across many environments. More importantly, we emphasize whole-body coordination by fully exploiting the dexterity of humanoid platforms. To this end, we propose the **Hu**manoid **M**anipulation **I**nterface (**HuMI**) (Fig. 1), a data-collection and learning framework with the following advantages:

**Robot-free, portable, and efficient data collection**: Our system requires only handheld sensorized grippers and base-station-free wearable pose trackers. This design enables task

---

†Work done during an internship at Tsinghua University.

teaching without the physical presence of a robot, and the entire setup fits into a single backpack. By eliminating the need to manage real-robot balance or manually compensate for controller tracking errors, our approach achieves a 3x increase in data-collection throughput compared to teleoperation [54].

**Broad task coverage and strong generalization**: HuMI supports a wide range of humanoid manipulation tasks involving whole-body coordination, precise bimanual actions, dynamic motions, and base mobility, requiring only changes in demonstration data. Furthermore, with diverse demonstrations collected across many environments, HuMI achieves a 70% success rate on unseen objects and environments.

Achieving these capabilities requires more than directly applying existing robot-free data collection systems [6, 12, 13] to humanoid robots. Traditional frameworks typically consist of (1) a data collection system that records end-effector trajectories, (2) a high-level policy that generates target trajectories from onboard observations, and (3) a low-level controller that executes these trajectories. However, humanoid whole-body manipulation introduces unique challenges that are not addressed by this pipeline. Below, we outline the key challenges and our strategies for addressing them.

**Underspecified demonstrations**: Existing robot-free data collection systems mainly target tasks involving one or two end-effectors (grippers). However, gripper trajectories alone are insufficient to specify whole-body manipulation. For example, squatting, kneeling, and bending can all achieve low-reaching motions, yet the movements of the waist, legs, and feet are often critical for success. To address this, we record trajectories not only for the grippers but also for the base (pelvis) and feet. We then use inverse kinematics (IK) to augment these trajectories into full robot degrees of freedom.

**Feasibility gap**: Morphological discrepancies often render human demonstrations kinematically infeasible, leading to issues such as self-collisions or reach limitations. Unlike traditional motion retargeting for expressive motions (e.g., dancing) [1, 26, 47], simply scaling the motion data is not viable for manipulation tasks, as the physical scene and object remain immutable. For example, scaling down arm length may result in the robot failing to reach a target. To ensure the kinematic feasibility of the original, unscaled trajectories, we develop an online IK preview interface that visualizes the resulting humanoid motion in real-time during data collection. This interface enables demonstrators to intuitively adjust their movements, ensuring the collected data is both task-compliant and executable.

**Non-negligible execution error**: Previous robot-free data collection frameworks rely on low-level controllers to execute human trajectories with high precision. However, despite advances in sim-to-real RL for humanoid trajectory tracking [2, 16, 23, 28, 37, 47, 54–56], non-negligible tracking errors (4–6 cm) persist. These errors compromise the original policy interfaces [6, 12, 13]. Specifically, high-level policies employing action chunking [7, 20, 57] exhibit discontinuities at chunk boundaries due to mismatches between planned and executed poses. To bridge this gap, we propose a manipulation-centric

whole-body controller designed to maximize precision without sacrificing stability, alongside a redesigned policy interface that improves the coordination between high and low-level controls.

We evaluate HuMI on five tasks: marriage proposal, squatting to pick up a bottle from the ground, tossing a toy, unsheathing a sword, and walking to clean a table. These tasks cover a wide range of whole-body manipulation behaviors. Our results demonstrate HuMI's high data-collection efficiency and strong task success rates. We further evaluate generalization and achieve a 70% success rate in unseen environments with unseen objects.

In summary, our contributions are:
- The first robot-free demonstration system for humanoid whole-body manipulation tasks.
- A learning framework enabling the transfer of manipulation skills from humans to humanoids by systematically overcoming the embodiment gap.
- Extensive real-world validation on five diverse whole-body tasks, demonstrating $3\times$ higher data-collection throughput compared to teleoperation and 70% success rates in unseen environments.

## II. METHOD

HuMI consists of two components: a robot-free demonstration system (Fig. 2), and a hierarchical policy learning framework (Fig. 3). First, we collect human demonstration data in the form of whole-body trajectories and image observations. These data are used to train the high-level manipulation policy, in which a Diffusion Policy [7] maps image observations to actions represented as target keypoint trajectories. The same data are also used to train the low-level controller, which outputs robot joint angles to track the target trajectories generated by the high-level policy. As shown in Fig. 3, by integrating the high-level policy with the low-level controller, the resulting system enables humanoid whole-body manipulation using observations from onboard sensors. In the following sections, we describe each component and their integration in detail.

### A. Robot-Free Demonstration System

The primary goal of our demonstration system is to capture informative and robot-feasible human trajectories without requiring the physical presence of a robot. To achieve this, the system integrates portable, precise task-space recording hardware with a data processing pipeline optimized for whole-body feasibility.

**Portable and precise hardware.** The hardware design of HuMI prioritizes portability and precision to capture raw data sufficiently rich for whole-body manipulation. We build upon UMI [6], a widely adopted robot-free data collection system using handheld grippers [8, 12, 13, 45]. However, relying solely on gripper trajectories is insufficient for specifying whole-body motion, as the configurations of the torso, waist, and legs are critical for task success. For instance, retrieving an object from under a table may require the robot to squat;
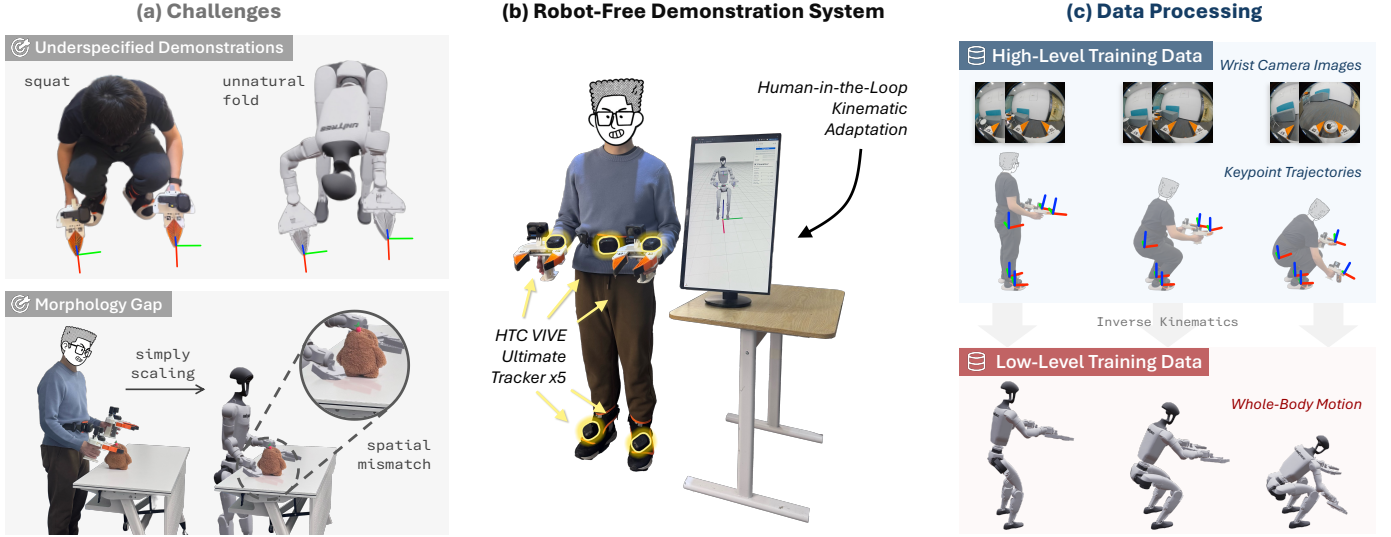
Fig. 2: **Overview of the HuMI data collection system.** (a) **Challenges**: Relying solely on gripper poses under-specifies whole-body motion, leading to unnatural postures (top); meanwhile, naively scaling human motions to match the robot's size compromises the spatial alignment required for object interaction (bottom). (b) **Hardware Setup**: Our portable system utilizes handheld sensorized grippers and trackers on the grippers, waist, and feet. A real-time IK preview interface enables human-in-the-loop kinematic adaptation. (c) **Data Processing**: Collected data serves two purposes: visual observations and task-space SE(3) trajectories train the high-level policy, while whole-body IK solutions provide reference motions for the low-level controller.

while bending might achieve the same end-effector pose (see Fig. 2 (a) upper), such postures appear unnatural and increase the risk of collision. Consequently, we adopt a standard full-body tracking configuration focusing on five key operational frames: the pelvis (floating base), hands, and feet[1] [3, 17, 44].

Unlike traditional outside-in Motion Capture systems [29, 43], which lack portability, our device must be standalone and base-station-free to enable data collection across diverse environments. Current standalone tracking solutions primarily categorize into headset-dependent systems (e.g., Pico [31]) and independent self-tracking systems (e.g., HTC Vive Ultimate Tracker [18]). We select the HTC Vive Ultimate Tracker to ensure robust whole-body tracking, as headset-based systems often suffer from tracking degradation during occlusion (e.g., when squatting to interact with ground-level objects). The resulting apparatus comprises two 3D-printed handheld grippers equipped with wrist-mounted GoPro cameras [6] and five trackers attached to the grippers, waist, and feet (Fig. 2 (b) and Fig. 8 in Appendix). As shown in Fig. 2 (c), the collected data includes synchronized image observations and task-space $SE(3)$ trajectories for grippers, base (pelvis), and feet, which drive the subsequent learning of the manipulation policy and low-level controller.

**Human-in-the-loop kinematic adaptation.** A core challenge for HuMI is overcoming the embodiment gap between the human operator and the humanoid robot. Traditional retargeting methods often scale human motions to match the robot's morphology [1, 26, 47]. However, scaling trajectories compromises the spatial relationship between the robot and the object, as the object's physical pose cannot be scaled. For instance, in Fig. 2 (a) bottom, simply scaling body heights and

arm length leads to insufficient reach and unintended intrusion. Although interaction geometry can theoretically be preserved using object meshes and poses [47], modeling and tracking every object is labor-intensive and costly. Furthermore, visuomotor whole-body manipulation requires strict visual-spatial alignment—visual perception must remain consistent with physical location. Therefore, HuMI focuses on tracking the original, unscaled poses from human trajectories.

Without scaling, however, these trajectories may become infeasible for the robot. Our target humanoid (Unitree G1 [41]) is approximately 130 cm tall; consequently, motions performed by an adult human may fall out of the robot's workspace, and self-collision risks increase when interacting with objects near the body. To ensure feasibility, we incorporate a human-in-the-loop adaptation mechanism via an online IK preview interface (see Fig. 2 (b)). By visualizing the virtual robot's kinematic motion in real-time, operators can adjust their demonstrations on the fly to satisfy both feasibility and task constraints. Unlike teleoperating a physical robot with complex dynamic constraints [2, 22, 52, 54], controlling a virtual robot subject only to kinematic constraints imposes a significantly lower cognitive load. This approach further benefits downstream learning by representing demonstrations with full-body degrees of freedom, providing comprehensive supervision for low-level controller training (see Fig. 2 (c)).

### B. Manipulation-Centric Whole-Body Controller

To execute target trajectories from the high-level policy, we train a reinforcement learning controller in simulation to track whole-body reference motions. Yet, state-of-the-art trackers [23, 28] often incur tracking deviations of 4–6 cm, which is insufficient for fine manipulation. While naively tightening end-effector (EE) tracking tolerance is intuitive, we find

---

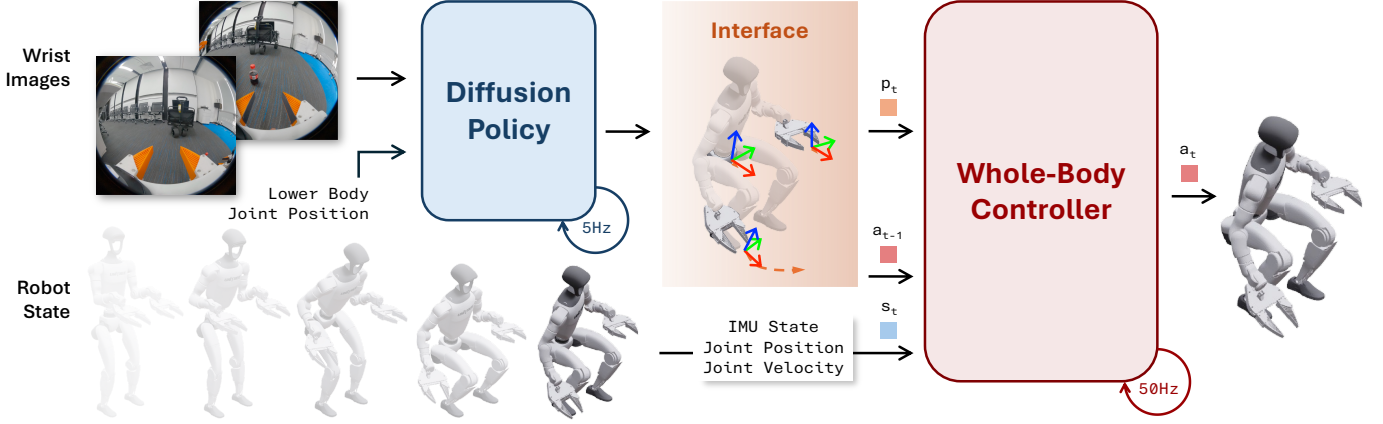[1]The head is excluded as the target robot [41] lacks an actuated neck.

Fig. 3: **Hierarchical control framework of HuMI.** (1) A high-level **Diffusion Policy** (5Hz) processes camera images and proprioception to generate receding-horizon task-space trajectories (action chunks). (2) A low-level **Whole-Body Controller** (50Hz) tracks these keypoint targets $p_t$, integrating the current robot state $s_t$ (IMU, joint positions/velocities) to compute precise joint actuation commands $a_t$.

it counterproductive: over-prioritizing end-effector precision leads to neglecting whole-body coordination, which actually compromises stability and impairs overall task performance (see Appendix C). To maintain coordination while seeking high precision, we introduce two mechanisms: adaptive tracking rewards and variable-speed augmentation.

**Adaptive end-effector tracking.** To learn coarse coordinated motion, we first employ a basic whole-body tracking reward $r(\bar{e}_\chi, \sigma_\chi) = \exp(-\bar{e}_\chi/\sigma_\chi^2)$ following standard practice. For each metric $\chi \in \{\mathbf{p}, \theta, v, w\}$—denoting position, orientation, and linear/angular velocity—$\bar{e}_\chi$ is the mean error defined as in [23] and the constant $\sigma_\chi$ denotes the precision tolerance. The total whole-body tracking reward is then:

$$r_{\text{body}} = \sum_{\chi \in \{\mathbf{p}, \theta, v, w\}} r(\bar{e}_\chi^{\text{body}}, \sigma_\chi).$$

Beyond basic coordination, manipulation tasks further require specialized precision for the end-effectors. Typically, these requirements often differ across motion phases. Consider a human kneeling to pick up an object: the initial rapid descent can be relatively loose, while the final grasp is slower to ensure a precise contact. Inspired by this intuition, we dynamically scale precision tolerance for the end-effectors: requiring high accuracy during slow interactions but granting greater flexibility as velocity increases. The adaptive end-effector reward is defined as:

$$r_{\text{EE}} = \mathbb{I}\left(\|v_{\text{base}}^{\text{ref}}\| < \delta\right) \cdot \sum_{\chi \in \{\mathbf{p}, \theta\}} r\left(\bar{e}_\chi^{\text{EE}}, \sigma_\chi\left(v_{\text{EE}}^{\text{ref}}\right)\right),$$

where the dynamic scaling term $\sigma_\chi(\hat{v}_{\text{EE}})$ is linearly interpolated between $[\sigma_\chi^{\min}, \sigma_\chi^{\max}]$ based on reference end-effector speed. This reward is further gated by $\mathbb{I}(\cdot)$, which deactivates end-effector tracking when the reference base velocity exceeds $\delta$ to prioritize stability during rapid movement. Combining the whole-body and end-effector objectives, the final tracking reward is formulated as:

$$r_{\text{tracking}} = w_{\text{body}} r_{\text{body}} + w_{\text{EE}} r_{\text{EE}}.$$

We also observed that a curriculum for the end-effector reward is necessary; otherwise, prematurely focusing on end-effector precision often leads to uncoordinated whole-body postures. Therefore, we gradually ramp up $w_{\text{EE}}$ and anneal $\sigma_\chi^{\min}$ during training, shifting the learning focus from stable global motion to precise EE alignment (see Appendix E for details).

**Variable-speed augmentation.** In standard motion tracking, the reference typically advances at a fixed speed. In this case, the target often moves on too fast before the policy can spend enough time fixing small mistakes, making it hard to learn highly precise movements. We therefore introduce a variable execution pace to overcome this limitation.

For a reference motion with duration $T$, we scale the execution speed within $[s_{\min}, s_{\max}]$ by sampling a new speed scaling factor $s_k$ every $\Delta$ seconds (see Appendix E for details). This variety of slow speeds within each episode gives the policy ample time to fix small errors, thereby facilitating the learning of high-precision movements.

### C. Policy Interface for Improved System Integration

As shown in Fig. 3, we implement the high-level policy using a Diffusion Policy [7] that predicts action chunks represented as relative keypoint trajectories [6, 12, 13]. However, we observe that naively feeding these targets to the low-level controller results in system fragility, where coupled errors from both levels can significantly compromise stability. To ensure robust whole-body execution, we introduce two critical modifications to the policy interface.

**Target pose as high-level action reference.** Even with improved tracking performance, the tracking error of the low-level controller remains non-negligible. A primary issue arising from this is action chunk discontinuity, as illustrated in Fig. 4. Previous manipulation policies typically use the actual EE pose as the reference frame for the current action chunk [6, 12, 13]. However, for whole-body humanoid manipulation, tracking errors create a discrepancy between the robot's *executed* pose (dark gray line) and the scheduled *target* pose (light gray line) at the chunk switching boundary ($t = 2$).
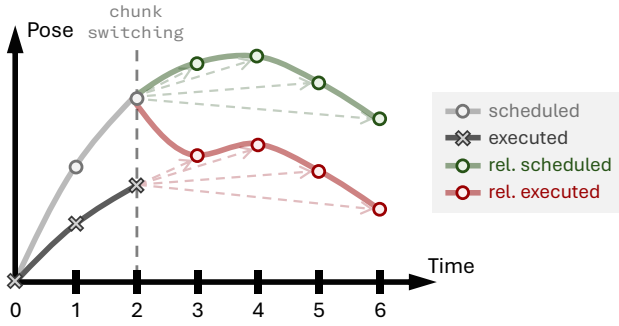
Fig. 4: **Impact of reference frame selection on action chunk continuity.** Due to tracking error, the executed robot pose (dark gray) "lags" behind the scheduled target (light gray). Naively anchoring the next action chunk to the current **executed** pose results in a sudden trajectory reversal (red line), disrupting momentum. By instead using the previous **scheduled** target as the reference frame, the policy produces a smooth, continuous trajectory (green line) that maintains the intended motion profile.

Consequently, resetting the reference to the lagging executed pose generates a trajectory that suffers from a sudden reversal (red line), disrupting the smooth momentum essential for dynamic tasks like tossing. To enforce continuity, we instead utilize the previous *target* pose as the action reference. As shown by the green line in Fig. 4, this approach naturally connects the current chunk with the previous one. Furthermore, this aligns better with the training dynamics of both levels: the high-level policy acts under the assumption of perfect tracking (as it is trained via imitation learning on human trajectories), while the low-level RL controller is trained to track fixed offline reference motions.
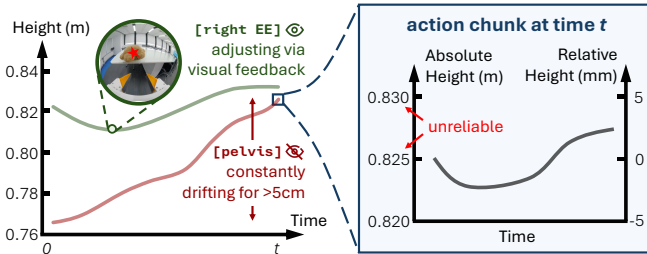


Fig. 5: **Mitigating drift in non-vision-grounded keypoints. Left**: Trajectories during a doll-grasping task. The "sighted" gripper (green) remains anchored via visual feedback, whereas the "blind" pelvis (red) suffers from open-loop drift ($> 5\,\mathrm{cm}$) over time. **Right**: Decomposition of the action chunk at time $t$. Because the absolute height (left axis) is corrupted by cumulative error, we discard absolute tracking in favor of relative transforms within the chunk (right axis).

**Relative pose tracking for non-vision-grounded keypoints.** Keypoint configuration is a critical design space for the policy interface. While previous frameworks often rely solely on gripper poses in the world frame [6, 12, 13], this is insufficient for whole-body manipulation. Theoretically, the HuMI demonstration data provides poses for full body keypoints; however, a trade-off exists between observability and the number of controlled keypoints. Unlike grippers equipped with wrist-view cameras, keypoints such as the pelvis and feet are "blind"—they lack direct visual anchors. Consequently, they accumulate unrecoverable errors during inference. As

illustrated in Fig. 5 (left), during a stationary grasping task, the target pelvis height drifts significantly ($> 5\,\mathrm{cm}$), rendering the absolute transform an unreliable control signal. To mitigate this, we modify the tracking objective for non-vision-grounded keypoints to track the relative transform within the current action chunk rather than the absolute transform (Fig. 5 right). This approach enables flexible keypoint configurations; we use a 3-keypoint setup (grippers + base) by default and demonstrate that the system maintains robust performance even when scaled to 5 keypoints (adding feet).

## III. EXPERIMENTS

In this section, we empirically evaluate HuMI along three key dimensions. Specifically, we aim to answer the following questions:

1) **Whole-body manipulation capability.** Can HuMI learn feasible whole-body skills from robot-free demonstrations, achieve sufficient manipulation precision while respecting whole-body dynamics, and effectively coordinate high-low level policy during fully autonomous execution?
2) **Generalization ability.** Do robot-free demonstrations collected across varied environments enable the learned policy to generalize to unseen environments and objects?
3) **Data-collection efficiency.** Can HuMI acquire whole-body manipulation data efficiently and with a high acceptance rate, and does the collected dataset cover versatile whole-body skills, including motions that are challenging to obtain for teleoperation?

To study these questions, we design five representative whole-body manipulation tasks and evaluate HuMI under both in-domain and out-of-domain settings with respect to environments and objects. To assess data-collection efficiency, we further compare HuMI against the state-of-the-art humanoid teleoperation system TWIST2[54].

## IV. WHOLE-BODY MANIPULATION CAPABILITY

In the capability experiments, we focus on evaluating HuMI's whole-body manipulation capability. We design four representative tasks that target complementary aspects—whole-body coordination, precise bimanual manipulation, high-speed dynamic motion, and long-range loco-manipulation. All experiments use in-domain settings, with tasks evaluated in the same environments and initial robot–object configurations as data collection.

### A. Learning Feasible Whole-Body Skills from Robot-Free Demonstrations

In the first experiment, we investigate whether HuMI can acquire feasible whole-body skills from robot-free demonstrations. We use a **marriage proposal** motion, in which the robot kneels from an upright stance onto its right knee, picks up a ring-shaped toy from the ground with its right hand, and raises it in a proposal gesture.

**Task challenges.** This task poses challenges for humanoid whole-body coordination. The robot must coordinate nearly
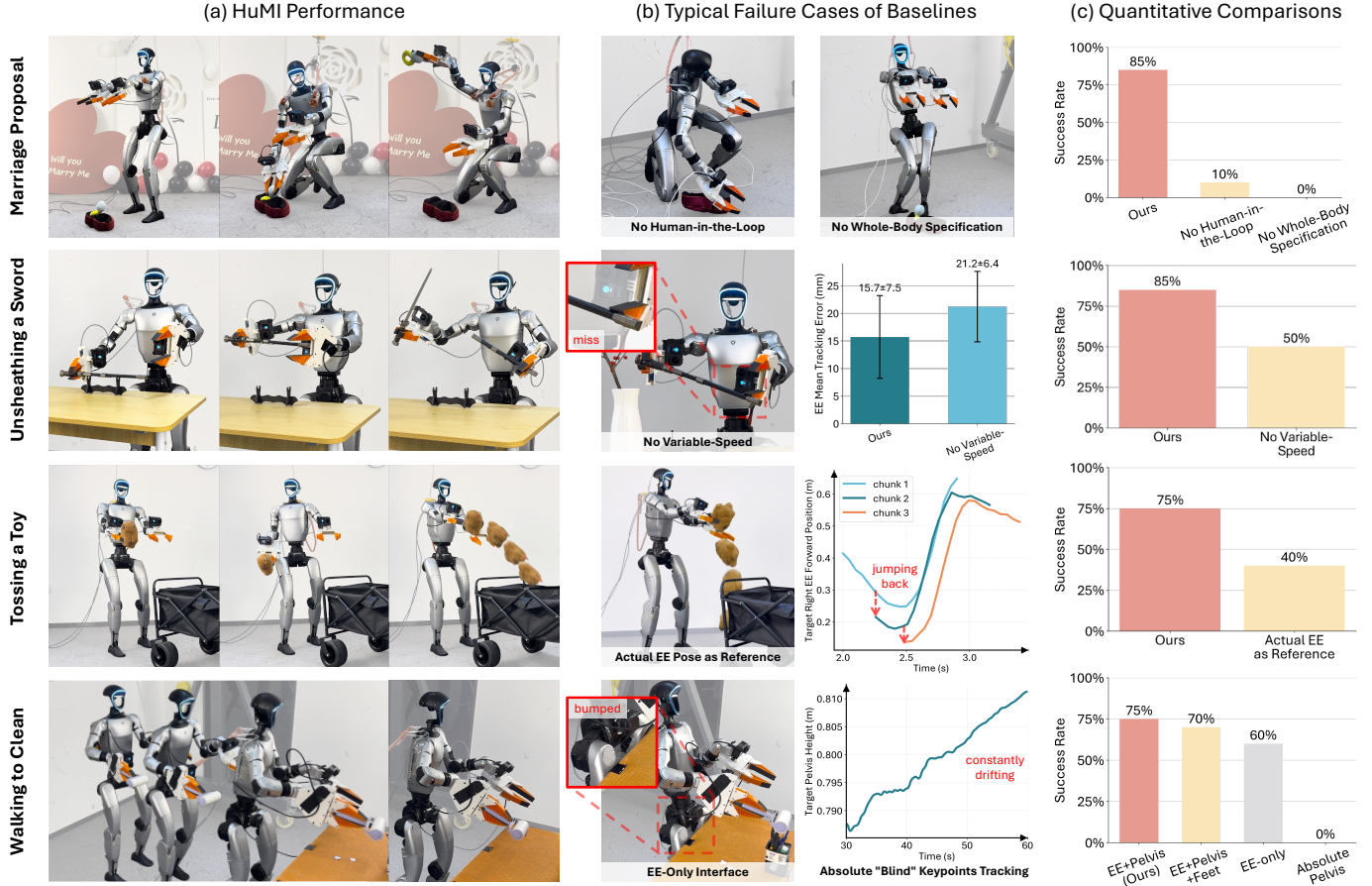
Fig. 6: **Whole-body capability experiment results**: (a) **HuMI performance** across the four tasks. (b) **Typical failure modes** of the ablation baselines. The red dashed box highlights the failure behavior details. (c) **Success rate** for each task.

all joints to transition from an upright stance to a single-knee kneel while keeping its center of mass within a very narrow support polygon and maintaining balance. At the lowest point, the robot must precisely grasp a small ring toy near the ground and lift it, requiring high end-effector accuracy under substantial whole-body movement.

**Performance.** HuMI successfully completes the marriage-proposal motion in $17/20 = 85\%$ cases (Fig. 6 (c)). The robot maintains balance and produces smooth, coordinated whole-body motion. The trajectories appear natural and human-like while still achieving a precise ring grasp and lift (Fig. 6 (a)).

**Remove human-in-the-loop kinematic adaptation.** To assess the role of online kinematic preview during data collection, we train HuMI on demonstrations recorded *without* the human-in-the-loop kinematic adaptation module. In this variant, operators no longer see the humanoid avatar and are thus not guided by the robot's reach or joint limits. Under this setting, the success rate drops from $17/20 = 85\%$ rollouts to $1/10 = 10\%$, and the learned policy often proposes kinematically inappropriate motions. As illustrated in Fig. 6 (b), the robot frequently kneels with an excessively splayed leg, highlighting the importance of kinematic feedback for keeping demonstrations within a humanoid-feasible and kinematic-appropriate motion space.

**Remove whole-body specification.** We further ablate re-

moving whole-body supervision. Following prior UMI-style interfaces [6, 12, 13], the high-level policy communicates only end-effector waypoints, and the low-level controller is trained from the two sparse EE waypoints. The success rate drops to $0/10$. The low-level controller struggles to maintain a stable, coordinated whole-body motion from these underspecified demonstrations and often converges to kinematically or dynamically inappropriate solutions.

### B. Precise Humanoid Bimanual Manipulation

To evaluate HuMI's capability for precise humanoid bimanual manipulation, we consider the **unsheathing a sword** motion. In this task, the robot first grasps the hilt on the rack, then stabilizes the scabbard in mid-air with its left hand and coordinates both arms to fully draw the blade.

**Task challenges.** This task demands accurate grasps, precise end-effector poses, and tightly coordinated bimanual motion within a narrow success basin. The sword hilt and scabbard offer small contact areas. Slight pose errors can destabilize contact or cause collisions with the rack. Once both hands are engaged, mismatches in timing or motion induce shear that degrades alignment, making this motion a stringent benchmark for precise humanoid bimanual manipulation.

**Performance.** HuMI successfully completes the unsheathing in $17/20 = 85\%$ trials, with an average end-effector tracking error of $15.7\,\text{mm}$ (Fig. 6 (c)). The robot secures

precise grasps on the hilt and scabbard without slippage or collisions with the rack. It then maintains tight bimanual coordination so that the blade slides out smoothly in a single continuous pull (Fig. 6 (a)).

**Remove variable-speed augmentation.** We ablate the variable-speed augmentation by training the low-level policy on demonstrations replayed only at original human speeds. The success rate drops from $85\%$ to $5/10 = 50\%$, with the average tracking error increasing to $21.2\,\text{mm}$ (Fig. 6 (b)). The policy often fails to secure a reliable grasp or lacks the precise bimanual synchronization. These degradations indicate that variable-speed augmentation is crucial for controller to resolve small tracking errors and improve bimanual coordination.

### C. Temporally Coherent Dynamic Control

The third experiment investigates whether HuMI can robustly and fluidly capture and transfer high-speed dynamic human motions to a humanoid robot. We consider **humanoid dynamic tossing** as the evaluation task. The robot must execute a temporally structured throwing motion, featuring a backward wind-up phase followed by a rapid forward swing that throws the object into a target container.

**Task challenges.** The robot must route momentum through many coupled joints, coordinating torso and arms for a backward wind-up and fast forward swing. These high-speed, continuous motions demand a tight control hierarchy: the high-level policy outputs smooth, temporally coherent trajectory commands, while the low-level policy fluidly realizes them as coordinated whole-body motion.

**Performance.** HuMI successfully completes the dynamic tossing task in $15/20 = 75\%$ trails (Fig. 6 (c)). The resulting throws exhibit smooth, stable whole-body trajectories, with the robot consistently releasing the object near the peak of the forward swing, producing accurate velocity and direction so that the object reliably lands inside the container (Fig. 6 (a)).

**Actual EE poses as action reference.** We ablate the action reference from target EE pose of the previous chunk to the actual executed EE pose. As shown in Fig. 6 (b)(c), the success rate drops from $15/20 = 75\%$ to $4/10 = 40\%$, and the resulting throws are noticeably more hesitant: direction reversals occur in the end-effector trajectory between chunks, which disrupt the monotonic forward-swing and introduce a jagged acceleration profile. Consequently, the object is often released with insufficient speed or a slightly incorrect direction, causing it to miss the container.

### D. Long-Range Loco-Manipulation

The last experiment evaluates the HuMI's long-range loco-manipulation capability. We consider **walking to clean the table** as a representative task. In each trial, the robot starts $1$–$2\,\text{m}$ away from the target desk, with its initial yaw offset sampled from $[-45°, 45°]$. The robot is required to navigate to the desk and then execute cleaning strokes with a lint roller to remove scattered paper scraps from the tabletop.

**Task challenges.** This task couples two distinct motion modalities: long-range walking and fine-grained tabletop

cleaning. The high-level policy first guides the robot toward the desk, then shifts its commands to focus on wiping once the surface is within reach. As this intent evolves, the keypoints interface must faithfully transmit it between the high-level and low-level policy, so the system can switch from prioritizing stable footsteps to precise wiping. This makes the approach-to-cleaning transition demanding, requiring final steps that leave the robot well positioned to sweep the tabletop.

**Performance.** Across 20 evaluation rollouts, HuMI successfully completes the loco-manipulation task in 15 cases (Fig. 6 (c)). In successful trials, the robot navigates from varied initial poses, positioning its final footsteps to ensure the tabletop remains well within the arm's workspace. It then performs overlapping wiping strokes with the lint roller, clearing all scattered paper scraps (Fig. 6 (a)).

**Design of keypoints interface.** We ablate the high–low-level interface across three paradigms: **EE-only**, which outputs end-effector targets; **EE+pelvis**, which additionally specifies a pelvis pose; and **EE+pelvis+feet**, which further adds feet targets. With the EE-only interface, the success rate drops to $6/10 = 60\%$ (Fig. 6 (b)(c)), with failures typically arising during the approach: the robot either stops too far from the desk or collides with it. With only end-effector targets, the low-level policy struggles to disambiguate locomotion-oriented commands (e.g., stepping forward) from manipulation-oriented commands (e.g., reaching forward), leading to inappropriate whole-body responses.

In contrast, the EE+pelvis and EE+pelvis+feet interfaces attain success rates of $15/20 = 75\%$ and $7/10 = 70\%$, respectively, with comparable behaviors. Providing additional body keypoints allows the high-level policy express richer whole-body intent and helps the low-level controller coordinate across motion modalities, so we treat both as viable interfaces within HuMI.

**Absolute pose tracking for non-vision-grounded keypoints.** We further probe our integration design by replacing relative pose tracking with absolute tracking for the non–vision-grounded keypoints. Using the EE+pelvis interface, we now track the pelvis in the global world frame instead of a relative frame. This change causes the success rate to collapse from $15/20 = 75\%$ with relative tracking to $0/10$. As shown in Fig. 6 (b), accumulated pelvis-tracking drift during the approach cannot be corrected without visual feedback, causing the robot to veer off course and ultimately lose balance.

## V. GENERALIZATION ABILITY

Existing humanoid loco-manipulation systems [2, 21, 52, 54] trained from teleoperated demonstrations typically collect data in a single, controlled lab. Consequently, training and evaluation often share near-identical environments and objects, leaving their generalization abilities unclear. In contrast, HuMI gathers robot-free demonstrations across diverse real-world scenes. We thus ask whether policies trained on such in-the-wild data can genuinely generalize to unseen configurations. We instantiate this study on a squat-and-pick-up task, where the humanoid visually localizes a floor-placed bottle and
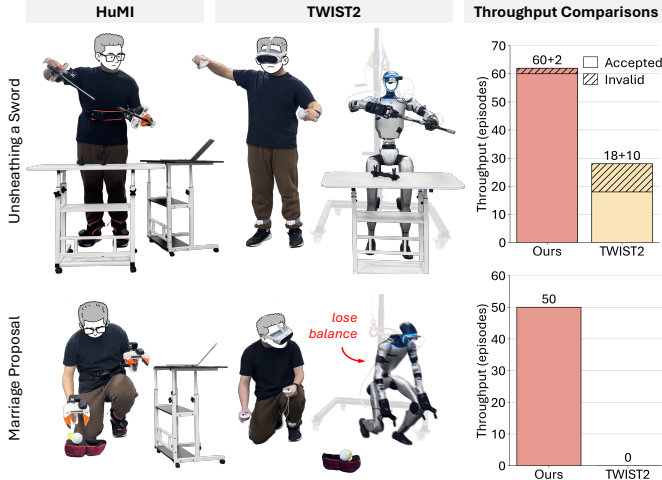
Fig. 7: **Data collection throughput comparison. Left**: HuMI and TWIST2 workflows. **Right**: Number of episodes collected within 15 min; dashed segments denote invalid trajectories.

guides its whole-body motion to squat, grasp, and lift it from near ground.

**Collecting whole-body manipulation demonstrations in various environments.** Thanks to the lightweight hardware and easy-to-deploy design of the HuMI data collection system, we can easily carry it into diverse real-world environments. We collect 350 whole-body demonstrations across 7 distinct environments and 7 different bottle instances, as shown in Appendix D. These in-the-wild demonstrations span variations in scene layout, lighting, and object appearance, and are used to train a diffusion policy for controlling the humanoid.

**Evaluation in generalization settings.** As shown in Appendix D, we evaluate the learned policy in two generalization regimes. (1) **Unseen environments.** We deploy the humanoid in four new scenes that differ from training locations in layout, clutter, and illumination, requiring the policy to extract task-relevant cues from camera observations despite these visual distractors. (2) **Unseen objects.** We further test on six novel items, including bottles and bottle-like objects absent from the training set. Across all trials, HuMI successfully completes $14/20 = 70\%$ episodes, maintaining reliable performance even in a dim stairwell and on out-of-distribution objects such as a vase whose shape and texture differ markedly from the training bottles.

## VI. DATA COLLECTION EFFICIENCY

In this section, we aim at quantifying how efficiently HuMI can capture demonstrations, how reliably users can complete recordings without failure, and to what extent the collected motions cover a broad spectrum of versatile whole-body manipulation behaviors. For comparison with conventional teleoperation pipelines, we use TWIST2[54] as a baseline, comparing its teleoperated workflow with HuMI in terms of efficiency, acceptance rate, and coverage.

**Throughput.** To evaluate data collection throughput, we use the unsheathing task as a shared benchmark and run 15 min collection sessions with both HuMI and TWIST2. For each system, we record the number of collected episodes and the

acceptance rate. To ensure a fair comparison, all sessions are conducted by experienced users (more than 20 h with HuMI and more than 10 h with TWIST2). We define the acceptance rate as the fraction of episodes that are usable for downstream humanoid policy learning: an episode is acceptable only if the trajectory successfully completes the task and a policy trained on the full set of collected trajectories can replay this trajectory end-to-end. As summarized in Fig. 7 (upper), HuMI yields substantially higher throughput, collecting **62** episodes versus **28** with TWIST2, while also achieving a higher acceptance rate (**96.7%** vs. **64.3%**). HuMI's streamlined, robot-free workflow further reduces the average time per acceptable episode to **30.0%** of that of TWIST2, indicating that users can obtain dense, high-quality datasets much more quickly with our robot-free pipeline.

**Whole-body motion coverage.** To evaluate the ability to capture versatile behaviors, we again use the marriage proposal motion as a challenging target task. As shown in Fig. 7 (bottom), HuMI successfully collected **50** demonstrations within 15 minutes with a **100%** acceptance rate, averaging just **18 s** per episode. Conversely, TWIST2 failed to produce any usable demonstrations, as the teleoperated humanoid cannot reliably realize the required deep kneeling and often lose stability. This underscores a key advantage of HuMI's robot-free data collection: it can capture diverse, highly articulated whole-body motions that go beyond the control limitations of existing humanoid teleoperation setups, providing broad coverage of complex behaviors for downstream policy learning.

## VII. RELATED WORKS

**Humanoid manipulation.** Prior research predominantly relies on sim-to-real RL [15, 25, 27, 46] or teleoperation [2, 4, 5, 9, 14, 21, 22, 52–54], yet both entail substantial overhead. Sim-to-real necessitates intricate reward shaping and domain randomization [30, 34, 40], while teleoperation imposes challenges for managing balance and compensating tracking errors. Although recent human video approaches [32, 48] show promise, they are largely restricted to hand motion transfer. In contrast, we propose a portable, robot-free system that enables robust whole-body manipulation across diverse tasks with strong generalization to unseen environments.

**Robot-free data collection.** This paradigm has shown high efficiency for fixed-base arms [6, 8, 19, 24, 35, 36, 38, 39, 45, 50] and recently floating-base platforms like quadrupeds [13] and aerial manipulators [12]. However, these methods typically only rely on end-effectors, lacking the capacity for complex whole-body coordination. We present the first robot-free data collection system specifically for humanoid whole-body manipulation.

## VIII. CONCLUSIONS AND LIMITATIONS

In this work, we introduced HuMI, a robot-free framework for data collection and learning in humanoid whole-body manipulation. Our system leverages portable hardware to capture whole-body motions without requiring the physical presence of a robot. By systematically addressing the embodiment

gap between humans and humanoids, our learning framework facilitates the transfer of diverse manipulation skills. We hope HuMI will contribute to democratizing humanoid data collection, improving learning efficiency, and fostering the development of more generalizable humanoid skills.

Despite its efficacy, limitations remain. First, the system relies on visual trackers [18], which require sufficient environmental texture and lighting. Second, while training configurations are unified, our low-level controllers are not yet general-purpose. Finally, evaluation was limited to a single platform [41]; however, we anticipate the framework can extend to other humanoids with minimal modification.

## IX. Acknowledgments

## References

[1] Joao Pedro Araujo, Yanjie Ze, Pei Xu, Jiajun Wu, and C. Karen Liu. Retargeting matters: General motion retargeting for humanoid motion tracking. *arXiv preprint arXiv:2510.02252*, 2025.

[2] Qingwei Ben, Feiyu Jia, Jia Zeng, Junting Dong, Dahua Lin, and Jiangmiao Pang. Homie: Humanoid loco-manipulation with isomorphic exoskeleton cockpit. *arXiv preprint arXiv:2502.13013*, 2025.

[3] Xiaowei Chen, Xiao Jiang, Lishuang Zhan, Shihui Guo, Qunsheng Ruan, Guoliang Luo, Minghong Liao, and Yipeng Qin. Full-body human motion reconstruction with sparse joint tracking using flexible sensors. *ACM Transactions on Multimedia Computing, Communications and Applications*, 20(2):1–19, 2023.

[4] Xuxin Cheng, Yandong Ji, Junming Chen, Ruihan Yang, Ge Yang, and Xiaolong Wang. Expressive whole-body control for humanoid robots. *arXiv preprint arXiv:2402.16796*, 2024.

[5] Xuxin Cheng, Jialong Li, Shiqi Yang, Ge Yang, and Xiaolong Wang. Open-television: Teleoperation with immersive active visual feedback. *arXiv preprint arXiv:2407.01512*, 2024.

[6] Cheng Chi, Zhenjia Xu, Chuer Pan, Eric Cousineau, Benjamin Burchfiel, Siyuan Feng, Russ Tedrake, and Shuran Song. Universal manipulation interface: In-the-wild robot teaching without in-the-wild robots. *arXiv preprint arXiv:2402.10329*, 2024.

[7] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, 44(10-11):1684–1704, 2025.

[8] Hojung Choi, Yifan Hou, Chuer Pan, Seongheon Hong, Austin Patel, Xiaomeng Xu, Mark R Cutkosky, and Shuran Song. In-the-wild compliant manipulation with umi-ft. *arXiv preprint arXiv:2601.09988*, 2026.

[9] Zipeng Fu, Qingqing Zhao, Qi Wu, Gordon Wetzstein, and Chelsea Finn. Humanplus: Humanoid shadowing and imitation from humans. *arXiv preprint arXiv:2406.10454*, 2024.

[10] Sergio Garrido-Jurado, Rafael Munoz-Salinas, Francisco José Madrid-Cuevas, and Rafael Medina-Carnicer. Generation of fiducial marker dictionaries using mixed integer linear programming. *Pattern recognition*, 51:481–491, 2016.

[11] GoPro, Inc. Hero10 black — waterproof action camera. https://gopro.com/en/us/shop/cameras/hero10-black/CHDHX-101-master.html, 2021. Accessed: 2026-01-29.

[12] Harsh Gupta, Xiaofeng Guo, Huy Ha, Chuer Pan, Muqing Cao, Dongjae Lee, Sebastian Scherer, Shuran Song, and Guanya Shi. Umi-on-air: Embodiment-aware guidance for embodiment-agnostic visuomotor policies. *arXiv preprint arXiv:2510.02614*, 2025.

[13] Huy Ha, Yihuai Gao, Zipeng Fu, Jie Tan, and Shuran Song. Umi on legs: Making manipulation policies mobile with manipulation-centric whole-body controllers. *arXiv preprint arXiv:2407.10353*, 2024.

[14] Tairan He, Zhengyi Luo, Xialin He, Wenli Xiao, Chong Zhang, Weinan Zhang, Kris Kitani, Changliu Liu, and Guanya Shi. Omnih2o: Universal and dexterous human-to-humanoid whole-body teleoperation and learning. *arXiv preprint arXiv:2406.08858*, 2024.

[15] Tairan He, Zi Wang, Haoru Xue, Qingwei Ben, Zhengyi Luo, Wenli Xiao, Ye Yuan, Xingye Da, Fernando Castañeda, Shankar Sastry, Changliu Liu, Guanya Shi, Linxi Fan, and Yuke Zhu. Viral: Visual sim-to-real at scale for humanoid loco-manipulation. *arXiv preprint arXiv:2511.15200*, 2025.

[16] Tairan He, Wenli Xiao, Toru Lin, Zhengyi Luo, Zhenjia Xu, Zhenyu Jiang, Jan Kautz, Changliu Liu, Guanya Shi, Xiaolong Wang, et al. Hover: Versatile neural whole-body controller for humanoid robots. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9989–9996. IEEE, 2025.

[17] Paul Heidicker, Eike Langbehn, and Frank Steinicke. Influence of avatar appearance on presence in social vr. In *2017 IEEE symposium on 3D user interfaces (3DUI)*, pages 233–234. IEEE, 2017.

[18] HTC VIVE. Vive ultimate tracker – full-body tracking for standalone vr, 2024. URL https://www.vive.com/eu/accessory/vive-ultimate-tracker/. Accessed: 2026-01-18.

[19] Yingdong Hu, Fanqi Lin, Pingyue Sheng, Chuan Wen, Jiacheng You, and Yang Gao. Data scaling laws in imitation learning for robotic manipulation. *arXiv preprint arXiv:2410.18647*, 2024.

[20] Lucy Lai, Ann Zixiang Huang, and Samuel J Gershman.

Action chunking as policy compression. *PsyArXiv*, 2022.

[21] Jialong Li, Xuxin Cheng, Tianshu Huang, Shiqi Yang, Ri-Zhao Qiu, and Xiaolong Wang. Amo: Adaptive motion optimization for hyper-dexterous humanoid whole-body control. *arXiv preprint arXiv:2505.03738*, 2025.

[22] Yixuan Li, Yutang Lin, Jieming Cui, Tengyu Liu, Wei Liang, Yixin Zhu, and Siyuan Huang. Clone: Closed-loop whole-body humanoid teleoperation for long-horizon tasks, 2025.

[23] Qiayuan Liao, Takara E Truong, Xiaoyu Huang, Yuman Gao, Guy Tevet, Koushil Sreenath, and C Karen Liu. Beyondmimic: From motion tracking to versatile humanoid control via guided diffusion. *arXiv preprint arXiv:2508.08241*, 2025.

[24] Fanqi Lin, Ruiqian Nai, Yingdong Hu, Jiacheng You, Junming Zhao, and Yang Gao. Onetwovla: A unified vision-language-action model with adaptive reasoning. *arXiv preprint arXiv:2505.11917*, 2025.

[25] Toru Lin, Kartik Sachdev, Linxi Fan, Jitendra Malik, and Yuke Zhu. Sim-to-real reinforcement learning for vision-based dexterous manipulation on humanoids. *arXiv:2502.20396*, 2025.

[26] Zhengyi Luo, Jinkun Cao, Alexander W. Winkler, Kris Kitani, and Weipeng Xu. Perpetual humanoid control for real-time simulated avatars. In *International Conference on Computer Vision (ICCV)*, 2023.

[27] Zhengyi Luo, Chen Tessler, Toru Lin, Ye Yuan, Tairan He, Wenli Xiao, Yunrong Guo, Gal Chechik, Kris Kitani, Linxi Fan, et al. Emergent active perception and dexterity of simulated humanoids from visual reinforcement learning. *arXiv preprint arXiv:2505.12278*, 2025.

[28] Zhengyi Luo, Ye Yuan, Tingwu Wang, Chenran Li, Sirui Chen, Fernando Castañeda, Zi-Ang Cao, Jiefeng Li, David Minor, Qingwei Ben, et al. Sonic: Supersizing motion tracking for natural humanoid whole-body control. *arXiv preprint arXiv:2511.07820*, 2025.

[29] OptiTrack. Primex 41 motion capture camera, 2024. URL https://optitrack.com/cameras/primex-41. Accessed: 2026-01-18.

[30] Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 3803–3810. IEEE, 2018.

[31] PICO. Pico virtual reality — official website, 2023. URL https://www.picoxr.com/global. Accessed: 2026-01-25.

[32] Ri-Zhao Qiu, Shiqi Yang, Xuxin Cheng, Chaitanya Chawla, Jialong Li, Tairan He, Ge Yan, David J. Yoon, Ryan Hoque, Lars Paulsen, Ge Yang, Jian Zhang, Sha Yi, Guanya Shi, and Xiaolong Wang. Humanoid policy ~ human policy. *arXiv preprint arXiv:2503.13441*, 2025.

[33] Francisco J Romero-Ramirez, Rafael Muñoz-Salinas, and Rafael Medina-Carnicer. Speeded up detection of squared fiducial markers. *Image and vision Computing*, 76:38–47, 2018.

[34] Fereshteh Sadeghi and Sergey Levine. Cad2rl: Real

single-image flight without a single real image. *arXiv preprint arXiv:1611.04201*, 2016.

[35] 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.

[36] Shuran Song, Andy Zeng, Johnny Lee, and Thomas Funkhouser. Grasping in the wild: Learning 6dof closed-loop grasping from low-cost demonstrations. *IEEE Robotics and Automation Letters*, 5(3):4978–4985, 2020.

[37] Zhi Su, Bike Zhang, Nima Rahmanian, Yuman Gao, Qiayuan Liao, Caitlin Regan, Koushil Sreenath, and S Shankar Sastry. Hitter: A humanoid table tennis robot via hierarchical planning and learning. *arXiv preprint arXiv:2508.21043*, 2025.

[38] Generalist AI Team. Gen-0: Embodied foundation models that scale with physical interaction. *Generalist AI Blog*, 2025. https://generalistai.com/blog/preview-uqlxvb-bb.html.

[39] RDT Team. Rdt2: Enabling zero-shot cross-embodiment generalization by scaling up umi data, September 2025. URL https://github.com/thu-ml/RDT2.

[40] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 23–30. IEEE, 2017.

[41] Unitree Robotics. Unitree G1: Humanoid Agent AI Avatar. https://www.unitree.com/g1, 2024. Accessed: 2026-01-16.

[42] Valve Corporation. Steamvr. https://store.steampowered.com/app/250820/SteamVR/, 2016. Accessed: 2026-01-29.

[43] Vicon Motion Systems. Valkyrie optical motion capture cameras, 2024. URL https://www.vicon.com/hardware/cameras/valkyrie/. Accessed: 2026-01-18.

[44] Timo Von Marcard, Bodo Rosenhahn, Michael J Black, and Gerard Pons-Moll. Sparse inertial poser: Automatic 3d human pose estimation from sparse imus. In *Computer graphics forum*, volume 36, pages 349–360. Wiley Online Library, 2017.

[45] Mengda Xu, Han Zhang, Yifan Hou, Zhenjia Xu, Linxi Fan, Manuela Veloso, and Shuran Song. Dexumi: Using human hand as the universal manipulation interface for dexterous manipulation. *arXiv preprint arXiv:2505.21864*, 2025.

[46] Haoru Xue, Tairan He, Zi Wang, Qingwei Ben, Wenli Xiao, Zhengyi Luo, Xingye Da, Fernando Castañeda, Guanya Shi, Shankar Sastry, et al. Opening the sim-to-real door for humanoid pixel-to-action policy transfer. *arXiv preprint arXiv:2512.01061*, 2025.

[47] Lujie Yang, Xiaoyu Huang, Zhen Wu, Angjoo Kanazawa, Pieter Abbeel, Carmelo Sferrazza, C Karen Liu, Rocky Duan, and Guanya Shi. Omniretarget: Interaction-preserving data generation for humanoid whole-body

loco-manipulation and scene interaction. *arXiv preprint arXiv:2509.26633*, 2025.

[48] Ruihan Yang, Qinxi Yu, Yecheng Wu, Rui Yan, Borui Li, An-Chieh Cheng, Xueyan Zou, Yunhao Fang, Xuxin Cheng, Ri-Zhao Qiu, et al. Egovla: Learning vision-language-action models from egocentric human videos. *arXiv preprint arXiv:2507.12440*, 2025.

[49] Brent Yi, Chung Min Kim, Justin Kerr, Gina Wu, Rebecca Feng, Anthony Zhang, Jonas Kulhanek, Hongsuk Choi, Yi Ma, Matthew Tancik, and Angjoo Kanazawa. Viser: Imperative, web-based 3d visualization in python. *arXiv preprint arXiv:2507.22885*, 2025.

[50] Sarah Young, Dhiraj Gandhi, Shubham Tulsiani, Abhinav Gupta, Pieter Abbeel, and Lerrel Pinto. Visual imitation made easy. In *Conference on Robot learning*, pages 1992–2005. PMLR, 2021.

[51] Kevin Zakka. Mink: Python inverse kinematics based on MuJoCo, December 2025. URL https://github.com/kevinzakka/mink.

[52] Yanjie Ze, Zixuan Chen, João Pedro Araújo, Zi ang Cao, Xue Bin Peng, Jiajun Wu, and C. Karen Liu. Twist: Teleoperated whole-body imitation system. *arXiv preprint arXiv:2505.02833*, 2025.

[53] Yanjie Ze, Zixuan Chen, Wenhao Wang, Tianyi Chen, Xialin He, Ying Yuan, Xue Bin Peng, and Jiajun Wu. Generalizable humanoid manipulation with 3d diffusion policies. In *2025 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2873–2880. IEEE, 2025.

[54] Yanjie Ze, Siheng Zhao, Weizhuo Wang, Angjoo Kanazawa, Rocky Duan, Pieter Abbeel, Guanya Shi, Jiajun Wu, and C Karen Liu. Twist2: Scalable, portable, and holistic humanoid data collection system. *arXiv preprint arXiv:2511.02832*, 2025.

[55] Tong Zhang, Boyuan Zheng, Ruiqian Nai, Yingdong Hu, Yen-Jen Wang, Geng Chen, Fanqi Lin, Jiongye Li, Chuye Hong, Koushil Sreenath, et al. Hub: Learning extreme humanoid balance. *arXiv preprint arXiv:2505.07294*, 2025.

[56] Siheng Zhao, Yanjie Ze, Yue Wang, C Karen Liu, Pieter Abbeel, Guanya Shi, and Rocky Duan. Resmimic: From general motion tracking to humanoid whole-body loco-manipulation via residual learning. *arXiv preprint arXiv:2510.05070*, 2025.

[57] 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 CONTENTS

## A. Data Collection System Details

This section provides details of our data collection system. We first describe the hardware components, followed by the human-in-the-loop IK adaptation interface, the data processing pipeline, and finally, the data collection protocol.



Fig. 8: **HuMI's hardware setup**.

*1) Hardware Components:* As illustrated in Fig. 8, our data collection system comprises the following components:

- **Two UMI [6] grippers**: We utilize UMI grippers equipped with GoPro cameras [11] to record wrist-view RGB observations and ArUco markers [10, 33] for gripper width tracking.
- **Five HTC VIVE Ultimate trackers [18]**: To capture 6-DoF poses, we attach trackers to the two grippers, the waist, and the feet. We modified the top cover of the original UMI gripper design to accommodate a tracker mount. Standard VIVE straps are used to secure the trackers to the waist and feet.

A laptop running SteamVR [42] is required to interface with the trackers to record their poses.

*2) IK Interface:* We developed a real-time IK preview interface to assist the demonstrator in adapting their motions to be kinematically feasible and task-compliant. Notably, to preserve strict spatial relationships between the robot and the environment, we use the original, unscaled motions (specifically, the trackers' recorded $SE(3)$ transforms in the world frame) as IK targets. Scaling is applied exclusively to the height of the pelvis tracker.

The IK problem incorporates three subtasks: tracking the translation and rotation of the five keypoints, avoiding self-collisions, and maintaining a natural configuration via posture
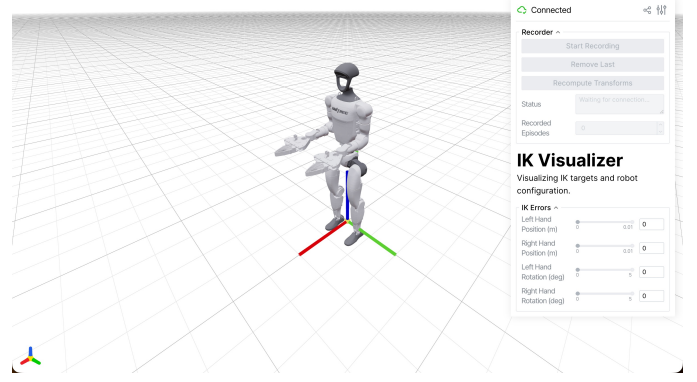


Fig. 9: **IK preview interface**.

regularization. We solve the IK problem in real-time using Mink [51] and visualize the robot's configuration using Viser [49]. Fig. 9 depicts the user interface during data collection.

*3) Data Processing:* The recorded data consists of MP4 videos from the gripper cameras and time-stamped $SE(3)$ trajectories from the five trackers. We process the data in the following steps:

1) **Synchronization**: We use the tracker timestamps as the reference clock since the five trackers are natively synchronized. To align the gripper videos with the tracker data, we extract the gyroscope data embedded in the video files by the GoPro cameras. We then align the video and tracker timestamps by cross-correlating the magnitudes of their angular velocities.
2) **Gripper width extraction**: Following Chi et al. [6], we extract the gripper width from the recorded videos by detecting the ArUco markers attached to the grippers.
3) **Data packaging**: The recorded data are packaged into two subsets: (1) visual observations, gripper widths, and keypoint trajectories for training the high-level policy; and (2) keypoint trajectories paired with whole-body IK solutions for training the low-level controller.

*4) Data Collection Protocol:* The following is the step-by-step protocol for collecting demonstrations in a new scene:

1) **Mapping setup**: For each new scene, the demonstrator must follow VIVE's prompts to build a tracking map of the environment; this typically takes 1–2 minutes.
2) **Calibration and synchronization**: The demonstrator may optionally perform gripper width calibration and GoPro timestamp synchronization following Chi et al. [6].
3) **Demonstration**: The demonstrator repeatedly performs the task within the scene. We use the start and stop times of the recorded videos to delineate episodes. The demonstrator uses GoPro's voice commands to control video recording, while the tracker data recording is controlled via the IK interface GUI (Fig. 9). Thanks to the synchronization step described in Sec. A3, strict alignment of start/stop times between the videos and tracker recordings is not required. During the demon-

stration, the demonstrator adapts their motions based on the real-time IK preview.

### B. Deployment Details

In this section, we detail the hardware infrastructure and software architecture employed in our real-world experiments. The experimental setup centers on the Unitree G1 humanoid robot [41], supported by an external workstation for high-level inference and HTC Vive Ultimate Trackers for global localization [18]. Below, we describe the custom end-effector design, the perception setup, and the hierarchical control system.

*1) Gripper Design:* To endow the Unitree G1 with manipulation capabilities while minimizing distal mass, we developed a custom gripper adapted from the UMI hardware interface [6]. We replaced the original spring-trigger mechanism with a direct-drive transmission inspired by the Wild LMA design. Specifically, we utilized the robot's existing wrist yaw motor to actuate the gripper, engineering a custom master gear that mates precisely with the motor's spline. This design allows us to remove the stock rubber hands and clamping mechanisms, securing the new mount via the original screw interfaces. While the transmission system was redesigned to enable direct actuation, the finger geometry and camera mount remain identical to the original UMI gripper to ensure compatibility with our training data.

*2) Camera Setup:* Visual observations are captured using two GoPro Hero 10 cameras [11] mounted on the grippers. Following the UMI hardware stack [6], we utilize a GoPro Media Mod to output HDMI signals, which are then converted to a low-latency USB 3.0 UVC interface via an Elgato HD60X capture card. These streams are transmitted to the workstation via USB, ensuring real-time observation updates.

*3) System Architecture:* As outlined in the main text, our control framework comprises a hierarchical structure: a high-level manipulation policy and a low-level whole-body controller. Communication between the external workstation and the robot's onboard computer is established via a local wireless network using ZeroMQ.

*a) High-Level Policy:* The high-level policy runs on the workstation at $5\,\mathrm{Hz}$. It aggregates visual streams from the external cameras and proprioceptive data (received from the robot) to infer desired end-effector keypoint trajectories and gripper commands. These targets are then published to the robot via the ZeroMQ interface.

*b) Low-Level Controller:* The low-level whole-body controller executes directly on the robot's onboard computer at $50\,\mathrm{Hz}$. Upon receiving the target keypoints and gripper commands, it computes the necessary joint position commands, which are executed by the robot's built-in PD controller. To support precise tracking, we attach one HTC Vive Ultimate Tracker to the robot's pelvis for global localization and place a second tracker on the ground to serve as a static $z = 0$ reference frame. Additional proprioceptive states, such as joint positions and IMU readings, are accessed directly from onboard sensors and streamed back to the high-level policy.

### C. Additional Experiments

In this section, we investigate the necessity of our adaptive end-effector (EE) reward. We aim to address a critical question raised in the methodology: whether naively prioritizing EE tracking precision in the current motion tracking framework would compromise whole-body stability. To validate this, we employ the **squat and pick up a bottle** task as a benchmark. This task is representative as it demands a seamless synergy between high-precision manipulation and whole-body dynamic balance: the humanoid must tightly coordinate its lower body and torso to maintain stability during the deep squat, while simultaneously achieving sufficient EE precision to execute a successful near-ground grasp and pick-up.

**Remove adaptive end-effector reward.** To instantiate the baseline, we establish a baseline that enforces a **naive tight tracking constraint** throughout the entire motion. Specifically, we replace the adaptive end-effector (EE) reward with a fixed formulation:

$$r_{\mathrm{EE}} = \sum_{\chi \in \{\mathbf{p}, \theta\}} r(\bar{e}_\chi^{\mathrm{EE}}, \sigma_\chi^{\mathrm{fixed}}),$$

where we set $\sigma_\chi^{\mathrm{fixed}}$ to the minimum tolerance used in our adaptive reward ($\sigma_{\mathbf{p}}^{\mathrm{fixed}} = 0.01\,\mathrm{m}$, $\sigma_\theta^{\mathrm{fixed}} = 5°$), and fix the weight to $\omega_{\mathrm{EE}}^{\mathrm{fixed}} = \omega_{\mathrm{EE}}^{\mathrm{max}} = 0.5$. This configuration enforces a uniformly tight EE tracking constraint throughout the motion, regardless of motion phase or modality. Under this setting, the success rate drops from $17/20 = 85\%$ to $5/10 = 50\%$. Typical failures occur during the deep-squat phase: the humanoid either loses balance and falls, or struggles to settle into an unnatural, marginally stable pose that prevents it from subsequently manipulation. This degradation highlights the critical role of the adaptive EE reward. By dynamically shaping the reward landscape according to motion modalities, our method balances the need of whole-body stability and manipulation precision throughout training. In contrast, permanently prioritizing EE precision forces the policy to neglect essential whole-body dynamics, compromising stability during whole-body dynamic phases like squatting.
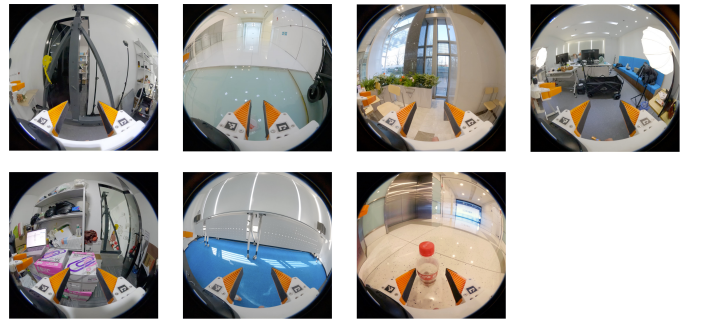
### D. Generalization Experiments Details



Fig. 10: **Training environments**.

**Environment details.** Fig. 10 visualizes the seven training environments used for our policy in Sec. V. We collected 50
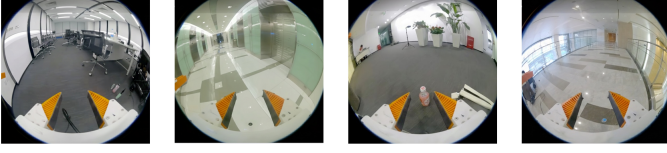
Fig. 11: **Testing environments**.

demonstrations in each environment, resulting in a total of 350 training trajectories. Fig. 11 visualizes the four testing environments, where we conducted 5 experiments in each environment.

**Object details.** Fig. 12 visualizes the seven training objects used for our policy in Sec. V. Fig. 13 visualizes the 6 testing objects.



Fig. 12: **Training objects**.



Fig. 13: **Testing objects**.

### E. Low-Level Controller Training Details

**Observation.** We train our low-level controller in a teacher–student framework. Similar to previous works [14, 16, 22, 52], we first train a teacher tracker that has access to privileged states and full-body reference commands. We then distill this teacher into a student policy that operates only on real-world states and keypoint trajectories aligned with the high-level policy, using the DAgger algorithm.

Specifically, at time step $t$, the teacher's observation is $O_t^{\text{tea}} = \left[s_t^{\text{tea}}, a_{t-1}^{\text{tea}}, c_t^{\text{tea}}\right]$, where the state $s_t^{\text{tea}} = [\mathbf{q}_t, \dot{\mathbf{q}}_t, \omega_t, g_t]$ includes the full-body joint positions and velocity $\mathbf{q}_t, \dot{\mathbf{q}}_t$, the base angular velocity $\omega_t$, and the base gravity vector $g_t$ projected into the body frame; $a_{t-1}^{\text{tea}}$ is the previous action. The whole-body command $c_t^{\text{tea}} = \left[\mathbf{q}_t^{\text{ref}}, \dot{\mathbf{q}}_t^{\text{ref}}, \mathbf{p}_t^{\text{ref}}, \theta_t^{\text{ref}}, \mathbf{p}_t^{\text{ref}} - \mathbf{p}_t, \theta_t^{\text{ref}} \ominus \theta_t\right]$ contains the reference joint positions and velocities $\mathbf{q}_t^{\text{ref}}, \dot{\mathbf{q}}_t^{\text{ref}}$, as well as the reference link positions and orientations $\mathbf{p}_t^{\text{ref}}, \theta_t^{\text{ref}}$ together with their deviations from the current link poses $\mathbf{p}_t, \theta_t$.

For the student policy, the observation is $O_t^{\text{stu}} = \left[s_{t-25:t}^{\text{stu}}, a_{t-26:t-1}^{\text{stu}}, c_t^{\text{stu}}\right]$, where $s_t^{\text{stu}} = [\mathbf{q}_t, v_t, \omega_t, g_t]$ denotes the same real-world state features as above, and we include a history of 25 steps of states and past actions to provide temporal context. The student command strictly aligns with the high-level policy output, $c_t^{\text{stu}} = \left[c_{\mathcal{T}_t}^{\text{EE}}, c_{\mathcal{T}_t}^{\text{blind}}\right]$, where each component contains a list of 10 waypoints sampled over the next 2 s:

$$\mathcal{T}_t = \{\, t_k = t + k \cdot \Delta t \mid k = 1, \ldots, 10 \,\}, \quad \Delta t = \left\lfloor \frac{2}{10\,\delta t} \right\rfloor,$$

and the control time step is $\delta t = 1/50\,\text{s}$.

The end-effector command at time step $t_k$ is $c_{t_k}^{\text{EE}} = \left[\mathbf{p}_{t_k}^{\text{ref}}, \theta_{t_k}^{\text{ref}}, \mathbf{p}_{t_k}^{\text{ref}} - \mathbf{p}_{t_k}, \theta_{t_k}^{\text{ref}} \ominus \theta_{t_k}\right]$, which includes the localized reference end-effector position and orientation $\mathbf{p}_{t_k}^{\text{ref}}, \theta_{t_k}^{\text{ref}}$, together with the position and orientation deltas, $\mathbf{p}_{t_k}^{\text{ref}} - \mathbf{p}_{t_k}$ and $\theta_{t_k}^{\text{ref}} \ominus \theta_{t_k}$. For "blind" points such as the pelvis or feet, we instead use relative displacements with respect to the current reference pose: $c_{t_k}^{\text{blind}} = \left[\mathbf{p}_{t_k}^{\text{ref}} - \mathbf{p}_t^{\text{ref}}, \theta_{t_k}^{\text{ref}} \ominus \theta_t^{\text{ref}}\right]$, which keeps the command *relative within each action chunk*, independent of the absolute position.

**Reward design.** Following prior works [23, 28], we decompose the low-level reward into a tracking term $r_{\text{tracking}}$ and a regularization term $r_{\text{penalty}}$. Table I summarizes all reward terms used for training the low-level controller. For the adaptive end-effector tracking rewards, we compute the tolerance parameter $\sigma_\chi$ by linearly interpolating between $\left[\sigma_\chi^{\min}, \sigma_\chi^{\max}\right]$ as a function of the commanded end-effector velocity $v_{\text{EE}}^{\text{ref}}$:

$$\sigma_\chi\left(v_{\text{EE}}^{\text{ref}}\right) = \text{clip}\left(f_{\text{interp}}(v_{\text{EE}}^{\text{ref}}), \sigma_\chi^{\min}, \sigma_\chi^{\max}\right),$$

$$f_{\text{interp}}(v_{\text{EE}}^{\text{ref}}) = \frac{v_{\text{EE}}^{\text{ref}} - v^{\min}}{v^{\max} - v^{\min}}(\sigma_\chi^{\max} - \sigma_\chi^{\min}) + \sigma_\chi^{\min}.$$

For the adaptive end-effector position tracking rewards, we set $\sigma_{\mathbf{p}}^{\min} = 0.01\,\text{m}$ and $\sigma_{\mathbf{p}}^{\max} = 0.1\,\text{m}$. For the adaptive end-effector rotation tracking rewards, we set $\sigma_\theta^{\min} = 5°$ and $\sigma_\theta^{\max} = 20°$. For both rewards, we set $v^{\max} = 0.1\,\text{m/s}$ and $v^{\min} = 0.05\,\text{m/s}$. Their reward weights are linearly increased from 0.0 to 0.5, and $\sigma_{\mathbf{p}}^{\min}$ is decreased from 0.1 m to 0.01 m over training steps 10,000 to 15,000.

**Domain randomization.** Table II summarizes the domain randomizations used for low-level controller training, grouped into physical, velocity, reset, and speed randomization.

The *reset pose shift* targets the mismatch between the command trajectory and the robot state at deployment. During training, the low-level policy tracks a replayed command that ignores the current state, while at test time the high-level policy issues online commands. When facing the temporal misalignment or sudden corrections, low-level controller may fail due to the out-of-distribution issue. To expose the policy to such cases, in each episode we sample a reset time $t_{\text{reset}}$ along the demonstration and initialize the robot to the pose at $t_{\text{reset}} + \mathcal{U}[-0.05, 0.05]$. This makes the robot start off from the reference yet still follow the same command trajectory.

For *variable-speed augmentation*, every 0.02 s we sample a speed scaling factor $s \sim \mathcal{N}_{[0.25, 1.25]}(1.0, \sigma^2)$ and use it to scale the reference time. The standard deviation $\sigma$ is linearly

TABLE I: **Rewards used in low-level controller training.**

| Reward Term | Equation | Weight |
|---|---|---|
| *Tracking Rewards* | | |
| Whole-body position tracking | $\exp(-\|\mathbf{p}^{\text{ref}} - \mathbf{p}\|_2^2/0.3^2)$ | 1.0 |
| Whole-body rotation tracking | $\exp(-\|\theta^{\text{ref}} \ominus \theta\|_2^2/0.4^2)$ | 1.0 |
| Whole-body linear velocity tracking | $\exp(-\|\mathbf{v}^{\text{ref}} - \mathbf{v}\|_2^2/1.0^2)$ | 1.0 |
| Whole-body angular velocity tracking | $\exp(-\|\omega^{\text{ref}} - \omega\|_2^2/\pi^2)$ | 1.0 |
| Adaptive end-effector position tracking | $\mathbb{I}(\|v_{\text{base}}^{\text{ref}}\| < 0.02\,\text{m/s}) \cdot \exp(-\|\mathbf{p}_{\text{EE}}^{\text{ref}} - \mathbf{p}_{\text{EE}}\|_2^2/\sigma_{\mathbf{p}}(v_{\text{EE}}^{\text{ref}})^2)$ | $0.0 \to 0.5$ |
| Adaptive end-effector rotation tracking | $\mathbb{I}(\|v_{\text{base}}^{\text{ref}}\| < 0.02\,\text{m/s}) \cdot \exp(-\|\theta_{\text{EE}}^{\text{ref}} \ominus \theta_{\text{EE}}\|_2^2/\sigma_{\theta}(v_{\text{EE}}^{\text{ref}})^2)$ | $0.0 \to 0.5$ |
| *Regularization Penalty* | | |
| Action rate | $\|a_t - a_{t-1}\|_2^2$ | $-5 \times 10^{-2}$ |
| Joint limits | $\sum \mathbb{I}(q_j \notin (q_j^{\min}, q_j^{\max}))$ | -10.0 |
| Undesired concact | $\sum_{i \notin \{\text{ankles,knees,hips}\}} \mathbb{I}(\|\mathbf{F}_i\|_2 > 1.0\,\text{N})$ | -0.1 |

TABLE II: **Domain randomization used in low-level controller training.**

| Term | Description | Sampling Range |
|---|---|---|
| *Physical randomization* | | |
| Static friction | randomize static friction of robot bodies | $\mu_s \sim \mathcal{U}[0.3, 1.6]$ |
| Dynamic friction | randomize dynamic friction of robot bodies | $\mu_d \sim \mathcal{U}[0.3, 1.2]$ |
| Resititution | randomize resitition friction of robot bodies | $e \sim \mathcal{U}[0.3, 1.2]$ |
| Default joint positions | add offsets to default joint positions | $\mathbf{q}_0 \sim \mathbf{q}_0 + \mathcal{U}[0.0, 0.5]$ |
| CoM offsets | randomize the torso link center of mass | $\Delta x \sim \mathcal{U}[-0.025, 0.025]$ $\Delta y \sim \mathcal{U}[-0.05, 0.05]$ $\Delta z \sim \mathcal{U}[-0.05, 0.05]$ |
| End-effector mass | randomize the mass of end-effector | $m \sim \mathcal{U}[0.75, 1.25] \cdot m$ |
| *Velocity perturbations* | | |
| Push robot | periodically push the robot every $\Delta t$ time | $\Delta t \sim \mathcal{U}[4, 6]$ $v_{\text{x,y}} \sim \mathcal{U}[-0.5, 0.5]$ $\omega_{\text{yaw}} \sim \mathcal{U}[-0.78, 0.78]$ |
| *Reset perturbations* | | |
| Base position | perturb base position at the reset time | $\Delta p_{\text{x,y}} \sim \mathcal{U}[-0.05, 0.05]$ $\Delta p_z \sim \mathcal{U}[0.0, 0.05]$ |
| Base orientation | perturb base orientation at the reset time | $\Delta\omega_{\text{pitch,roll}} \sim \mathcal{U}[-0.1, 0.1]$ $\Delta\omega_{\text{yaw}} \sim \mathcal{U}[-0.2, 0.2]$ |
| Base linear velocity | perturb base linear velocity at the reset time | $\Delta v_{\text{x,y}} \sim \mathcal{U}[-0.05, 0.05]$ $\Delta v_z \sim \mathcal{U}[-0.2, 0.2]$ |
| Base rotation velocity | perturb base rotation velocity at the reset time | $\Delta\omega_{\text{roll,pitch}} \sim \mathcal{U}[-0.52, 0.52]$ $\Delta\omega_{\text{yaw}} \sim \mathcal{U}[-0.78, 0.78]$ |
| Joint position | perturb joint positions at the reset time | $\Delta q \sim \mathcal{U}[-0.1, 0.1]$ |
| Reset pose shift | shift reset pose along the trajectory around $t_{\text{reset}}$ | $t_{\text{reset}} \sim t_{\text{reset}} + \mathcal{U}[-0.05, 0.05]$ |
| *Speed randomization* | | |
| Variable-speed augmentation | sample a play speed every 0.02s | $s \sim \mathcal{N}_{[0.25, 1.25]}(\mu = 1.0, \ \sigma^2)$ |

increased from $10^{-4}$ to 1.0 between training steps 10,000 and 15,000, after which it is kept at $\sigma = 1.0$.

*F. High-Level Policy Training Details*

We employ Diffusion Policy [7] as the backbone for high-level manipulation. Notably, this component is designed with modularity in mind; alternative architectures, such as ACT [57], can be seamlessly substituted.

*a) Hyperparameters:* The training hyperparameters remain consistent across all experimental tasks and are summarized in Table III.

*b) Observation and Action Spaces:* The high-level policy conditions on both visual and proprioceptive observations. Visual inputs comprise RGB images captured by two gripper-mounted cameras, each with a resolution of $224 \times 224$ pixels. Proprioceptive observations consist of the robot's lower-body joint angles. The action space defines the desired positions and rotations of the keypoints, along with gripper widths. We adopt the same action parameterization as UMI [6].

TABLE III: **Hyperparameters for the high-level diffusion policy.**

| Hyperparameter | Value |
|---|---|
| Visual observation horizon | 1 |
| Visual observation frequency | 20 Hz |
| Proprioceptive observation horizon | 3 |
| Proprioceptive observation frequency | 20 Hz |
| Action horizon | 48 |
| Action frequency | 20 Hz |
| Execution-to-data speed ratio | $1\times$ |
| Image resolution ($N_{\text{cam}} \times H \times W$) | $2 \times 224 \times 224$ |
| Vision backbone | `vit_base_patch14_dinov2.lvd142m` |
| Diffusion Policy learning rate | 3e-4 |
| Vision backbone learning rate | 3e-5 |
| Epochs | 200 |
| Batch size | 256 |
| Training loss | Flow Matching |
| Inference denoising steps | 10 |

*c) Training Data:* We utilize 100 demonstrations collected in a single environment for the capability tasks. These tasks include marriage proposal, unsheathing a sword, tossing a toy, and walking to clean a table. For the task of squatting to pick up a bottle, we use 350 demonstrations collected across seven distinct environments ($7 \times 50$).