

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

Anonymous CVPR submission

Paper ID 2

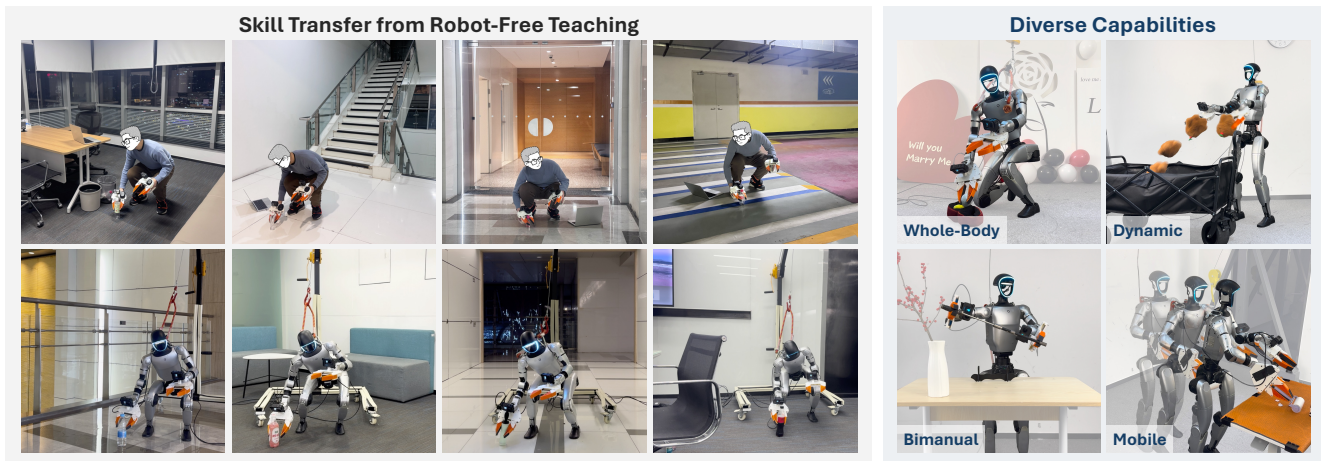


Figure 1. **Humanoid Manipulation Interface (HuMI)**. **Left:** Our portable, robot-free data collection facilitates skill transfer from human to humanoid across diverse, unstructured environments. **Right:** The framework enables a wide repertoire of complex whole-body behaviors.

Abstract

001 *Current approaches for humanoid whole-body manipula-*
 002 *tion, primarily relying on teleoperation or visual sim-to-*
 003 *real reinforcement learning, are hindered by hardware lo-*
 004 *gistics and complex reward engineering. Consequently,*
 005 *demonstrated autonomous skills remain limited and are*
 006 *typically restricted to controlled environments. In this*
 007 *paper, we present the Humanoid Manipulation Interface*
 008 *(HuMI), a portable and efficient framework for learning*
 009 *diverse whole-body manipulation tasks across various en-*
 010 *vironments. HuMI enables robot-free data collection by*
 011 *capturing rich whole-body motion using portable hard-*
 012 *ware. This data drives a hierarchical learning pipeline*
 013 *that translates human motions into dexterous and feasible*
 014 *humanoid skills. Extensive experiments across five whole-*
 015 *body tasks—including kneeling, squatting, tossing, walk-*
 016 *ing, and bimanual manipulation—demonstrate that HuMI*
 017 *achieves a 3x increase in data collection efficiency com-*
 018 *pared to teleoperation and attains a 70% success rate in*
 019 *unseen environments.*

1. Introduction

Humans expertly coordinate their entire bodies for manip-
 ulation, whether squatting to retrieve objects or bending to
 reach low tables. With their high degrees of freedom, hu-
 manoid 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 imita-
 tion learning from teleoperation [2, 21, 22, 52, 54]. How-
 ever, these methods are labor-intensive: RL demands di-
 verse assets and meticulous reward engineering, while tele-
 operation requires significant expertise to manage balance
 and controller inaccuracies. Consequently, current meth-
 ods demonstrate few autonomous tasks in fixed lab envi-
 ronments [2, 15, 21, 25, 46, 54]. These tasks exhibit limited
 whole-body coordination, typically restricting robots to up-
 right walking combined with simple actions like transport-
 ing objects, opening doors, or kicking boxes.

In this project, our goal is to enable humanoid robots to
 perform diverse tasks across many environments. More im-
 portantly, we emphasize whole-body coordination by fully

042 exploiting the dexterity of humanoid platforms. To this end,
043 we propose the **Humanoid Manipulation Interface (HuMI)**
044 (Fig. 1), a data-collection and learning framework with the
045 following advantages:

046 **Robot-free, portable, and efficient data collection:** Our
047 system requires only handheld sensorized grippers and
048 base-station-free wearable pose trackers. This design en-
049 ables task teaching without the physical presence of a robot,
050 and the entire setup fits into a single backpack. By elim-
051 inating the need to manage real-robot balance or manu-
052 ally compensate for controller tracking errors, our approach
053 achieves a 3x increase in data-collection throughput com-
054 pared to teleoperation [54].

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

063 Achieving these capabilities requires more than directly
064 applying existing robot-free data collection systems [6, 12,
065 13] to humanoid robots. Traditional frameworks typi-
066 cally consist of (1) a data collection system that records
067 end-effector trajectories, (2) a high-level policy that gen-
068 erates target trajectories from onboard observations, and
069 (3) a low-level controller that executes these trajectories.
070 However, humanoid whole-body manipulation introduces
071 unique challenges that are not addressed by this pipeline.
072 Below, we outline the key challenges and our strategies for
073 addressing them.

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

085 **Feasibility gap:** Morphological discrepancies often render
086 human demonstrations kinematically infeasible, leading to
087 issues such as self-collisions or reach limitations. Unlike
088 traditional motion retargeting for expressive motions (e.g.,
089 dancing) [1, 26, 47], simply scaling the motion data is not
090 viable for manipulation tasks, as the physical scene and ob-
091 ject remain immutable. For example, scaling down arm
092 length may result in the robot failing to reach a target. To
093 ensure the kinematic feasibility of the original, unscaled tra-
094 jectories, we develop an online IK preview interface that vi-

095 sualizes the resulting humanoid motion in real-time during
096 data collection. This interface enables demonstrators to in-
097 tuitively adjust their movements, ensuring the collected data
098 is both task-compliant and executable.

099 **Non-negligible execution error:** Previous robot-free data
100 collection frameworks rely on low-level controllers to exe-
101 cute human trajectories with high precision. However, de-
102 spite advances in sim-to-real RL for humanoid trajectory
103 tracking [2, 16, 23, 28, 37, 47, 54–56], non-negligible track-
104 ing errors (4–6 cm) persist. These errors compromise the
105 original policy interfaces [6, 12, 13]. Specifically, high-
106 level policies employing action chunking [7, 20, 57] exhibit
107 discontinuities at chunk boundaries due to mismatches be-
108 tween planned and executed poses. To bridge this gap, we
109 propose a manipulation-centric whole-body controller de-
110 signed to maximize precision without sacrificing stability,
111 alongside a redesigned policy interface that improves the
112 coordination between high and low-level controls.

113 We evaluate HuMI on five tasks: marriage proposal,
114 squatting to pick up a bottle from the ground, tossing a
115 toy, unsheathing a sword, and walking to clean a table.
116 These tasks cover a wide range of whole-body manipula-
117 tion behaviors. Our results demonstrate HuMI’s high data-
118 collection efficiency and strong task success rates. We fur-
119 ther evaluate generalization and achieve a 70% success rate
120 in unseen environments with unseen objects.

121 In summary, our contributions are:

- 122 • The first robot-free demonstration system for humanoid
123 whole-body manipulation tasks.
- 124 • A learning framework enabling the transfer of manipula-
125 tion skills from humans to humanoids by systematically
126 overcoming the embodiment gap.
- 127 • Extensive real-world validation on five diverse whole-
128 body tasks, demonstrating $3\times$ higher data-collection
129 throughput compared to teleoperation and 70% success
130 rates in unseen environments.

131 2. Method

132 HuMI consists of two components: a robot-free demon-
133 stration system (Fig. 2), and a hierarchical policy learning
134 framework (Fig. 3). First, we collect human demonstra-
135 tion data in the form of whole-body trajectories and image
136 observations. These data are used to train the high-level
137 manipulation policy, in which a Diffusion Policy [7] maps
138 image observations to actions represented as target keypoint
139 trajectories. The same data are also used to train the low-
140 level controller, which outputs robot joint angles to track
141 the target trajectories generated by the high-level policy. As
142 shown in Fig. 3, by integrating the high-level policy with the
143 low-level controller, the resulting system enables humanoid
144 whole-body manipulation using observations from onboard
145 sensors. In the following sections, we describe each com-
146 ponent and their integration in detail.

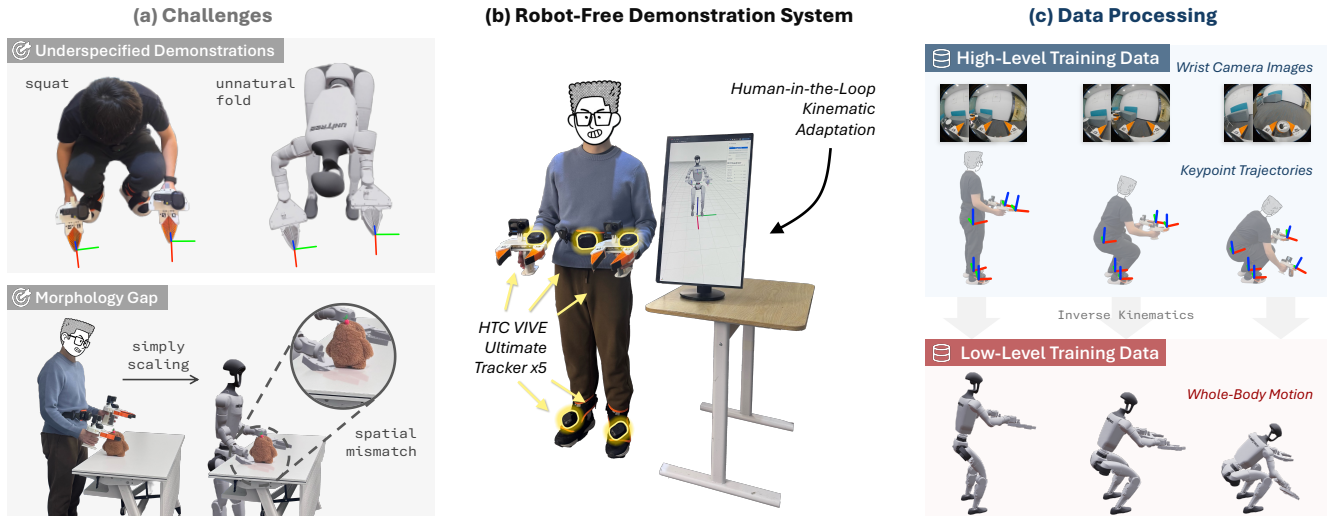


Figure 2. **Overview of the HuMI data collection system.** (a) **Challenges:** Relying solely on gripper poses under-specifies whole-body motion, leading to unnatural postures (top); meanwhile, naively scaling human motions to match the robot’s size compromises the spatial alignment required for object 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.

147 2.1. Robot-Free Demonstration System

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

154 **Portable and precise hardware.** The hardware design
155 of HuMI prioritizes portability and precision to capture
156 raw data sufficiently rich for whole-body manipulation.
157 We build upon UMI [6], a widely adopted robot-free data
158 collection system using handheld grippers [8, 12, 13, 45].
159 However, relying solely on gripper trajectories is insuffi-
160 cient for specifying whole-body motion, as the configura-
161 tions of the torso, waist, and legs are critical for task suc-
162 cess. For instance, retrieving an object from under a table
163 may require the robot to squat; while bending might
164 achieve the same end-effector pose (see Fig. 2 (a) upper),
165 such postures appear unnatural and increase the risk of colli-
166 sion. Consequently, we adopt a standard full-body tracking
167 configuration focusing on five key operational frames: the
168 pelvis (floating base), hands, and feet¹ [3, 17, 44].

169 Unlike traditional outside-in Motion Capture systems
170 [29, 43], which lack portability, our device must be stan-
171 dalone and base-station-free to enable data collection across
172 diverse environments. Current standalone tracking solu-
173 tions primarily categorize into headset-dependent systems
174 (e.g., Pico [31]) and independent self-tracking systems
175 (e.g., HTC Vive Ultimate Tracker [18]). We select the HTC

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

176 Vive Ultimate Tracker to ensure robust whole-body track-
177 ing, as headset-based systems often suffer from tracking
178 degradation during occlusion (e.g., when squatting to in-
179 teract with ground-level objects). The resulting apparatus
180 comprises two 3D-printed handheld grippers equipped with
181 wrist-mounted GoPro cameras [6] and five trackers attached
182 to the grippers, waist, and feet (Fig. 2 (b) and Fig. 8 in Ap-
183 pendix). As shown in Fig. 2 (c), the collected data includes
184 synchronized image observations and task-space $SE(3)$ tra-
185 jectories for grippers, base (pelvis), and feet, which drive
186 the subsequent learning of the manipulation policy and low-
187 level controller.

188 **Human-in-the-loop kinematic adaptation.** A core
189 challenge for HuMI is overcoming the embodiment gap be-
190 tween the human operator and the humanoid robot. Tra-
191 ditional retargeting methods often scale human motions to
192 match the robot’s morphology [1, 26, 47]. However, scaling
193 trajectories compromises the spatial relationship between
194 the robot and the object, as the object’s physical pose can-
195 not be scaled. For instance, in Fig. 2 (a) bottom, sim-
196 ply scaling body heights and arm length leads to insuffi-
197 cient reach and unintended intrusion. Although interac-
198 tion geometry can theoretically be preserved using object
199 meshes and poses [47], modeling and tracking every ob-
200 ject is labor-intensive and costly. Furthermore, visuomo-
201 tor whole-body manipulation requires strict visual-spatial
202 alignment—visual perception must remain consistent with
203 physical location. Therefore, HuMI focuses on tracking the
204 original, unscaled poses from human trajectories.

Without scaling, however, these trajectories may become

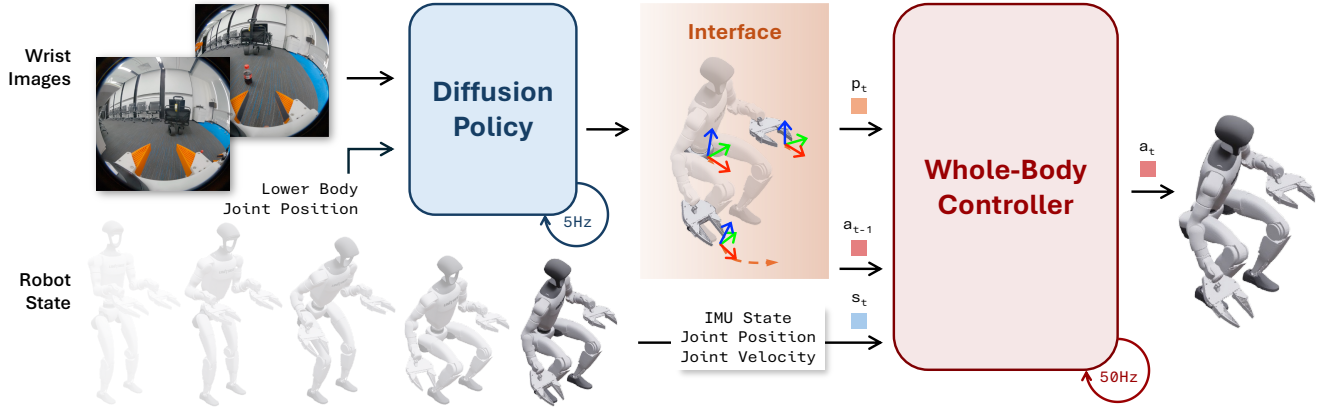


Figure 3. **Hierarchical control framework of HuMI.** (1) A high-level **Diffusion Policy** (5Hz) processes camera images and proprioception to generate receding-horizon task-space trajectories (action chunks). (2) A low-level **Whole-Body Controller** (50Hz) tracks these keypoint targets p_t , integrating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

206 infeasible for the robot. Our target humanoid (Unitree G1
207 [41]) is approximately 130 cm tall; consequently, motions
208 performed by an adult human may fall out of the robot’s
209 workspace, and self-collision risks increase when interact-
210 ing with objects near the body. To ensure feasibility, we
211 incorporate a human-in-the-loop adaptation mechanism via
212 an online IK preview interface (see Fig. 2 (b)). By visu-
213 alizing the virtual robot’s kinematic motion in real-time,
214 operators can adjust their demonstrations on the fly to sat-
215 isfy both feasibility and task constraints. Unlike teleoper-
216 ating a physical robot with complex dynamic constraints
217 [2, 22, 52, 54], controlling a virtual robot subject only to
218 kinematic constraints imposes a significantly lower cogni-
219 tive load. This approach further benefits downstream learn-
220 ing by representing demonstrations with full-body degrees
221 of freedom, providing comprehensive supervision for low-
222 level controller training (see Fig. 2 (c)).

2.2. Manipulation-Centric Whole-Body Controller

224 To execute target trajectories from the high-level policy,
225 we train a reinforcement learning controller in simulation
226 to track whole-body reference motions. Yet, state-of-the-
227 art trackers [23, 28] often incur tracking deviations of 4–
228 6 cm, which is insufficient for fine manipulation. While
229 naively tightening end-effector (EE) tracking tolerance is
230 intuitive, we find it counterproductive: over-prioritizing
231 end-effector precision leads to neglecting whole-body coor-
232 dination, which actually compromises stability and impairs
233 overall task performance (see Appendix C). To maintain co-
234 ordination while seeking high precision, we introduce two
235 mechanisms: adaptive tracking rewards and variable-speed
236 augmentation.

Adaptive end-effector tracking. To learn coarse coordi-
nated motion, we first employ a basic whole-body track-
ing 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 po-
sition, orientation, and linear/angular velocity— \bar{e}_χ is the
mean error defined as in [23] and the constant σ_χ denotes
the precision tolerance. The total whole-body tracking re-
ward 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. Typi-
cally, 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 intu-
ition, 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:

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

where the dynamic scaling term $\sigma_\chi(v_{\text{EE}})$ is linearly inter-
polated between $[\sigma_\chi^{\text{min}}, \sigma_\chi^{\text{max}}]$ based on reference end-effector
speed. This reward is further gated by $\mathbb{I}(\cdot)$, which deacti-
vates end-effector tracking when the reference base veloc-
ity exceeds δ 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_{EE} and

241 anneal σ_{χ}^{\min} during training, shifting the learning focus
 242 from stable global motion to precise EE alignment (see Ap-
 243 pendix E for details).

244 **Variable-speed augmentation.** In standard motion
 245 tracking, the reference typically advances at a fixed speed.
 246 In this case, the target often moves on too fast before the
 247 policy can spend enough time fixing small mistakes, mak-
 248 ing it hard to learn highly precise movements. We therefore
 249 introduce a variable execution pace to overcome this limita-
 250 tion.

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

257 2.3. Policy Interface for Improved System Integra- 258 tion

259 As shown in Fig. 3, we implement the high-level policy using
 260 a Diffusion Policy [7] that predicts action chunks repre-
 261 sented as relative keypoint trajectories [6, 12, 13]. How-
 262 ever, we observe that naively feeding these targets to the
 263 low-level controller results in system fragility, where cou-
 264 pled errors from both levels can significantly compromise
 265 stability. To ensure robust whole-body execution, we intro-
 266 duce two critical modifications to the policy interface.

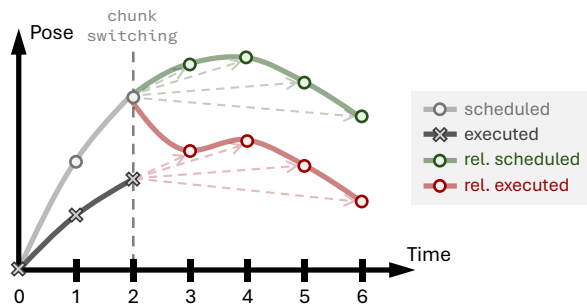


Figure 4. **Impact of reference frame selection on action chunk continuity.** Due to tracking error, the executed robot pose (dark gray) “lags” behind the scheduled target (light gray). Naively anchoring the next action chunk to the current **executed** pose results in a sudden trajectory reversal (red 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.

267 **Target pose as high-level action reference.** Even with
 268 improved tracking performance, the tracking error of the
 269 low-level controller remains non-negligible. A primary is-
 270 sue arising from this is action chunk discontinuity, as illus-
 271 trated in Fig. 4. Previous manipulation policies typically
 272 use the actual EE pose as the reference frame for the cur-

273 rent action chunk [6, 12, 13]. However, for whole-body hu-
 274 manoid manipulation, tracking errors create a discrepancy
 275 between the robot’s *executed* pose (dark gray line) and the
 276 scheduled *target* pose (light gray line) at the chunk switch-
 277 ing boundary ($t = 2$). Consequently, resetting the reference
 278 to the lagging executed pose generates a trajectory that suf-
 279 fers from a sudden reversal (red line), disrupting the smooth
 280 momentum essential for dynamic tasks like tossing. To en-
 281 force continuity, we instead utilize the previous *target* pose
 282 as the action reference. As shown by the green line in Fig. 4,
 283 this approach naturally connects the current chunk with the
 284 previous one. Furthermore, this aligns better with the train-
 285 ing dynamics of both levels: the high-level policy acts under
 286 the assumption of perfect tracking (as it is trained via imita-
 287 tion learning on human trajectories), while the low-level RL
 288 controller is trained to track fixed offline reference motions.

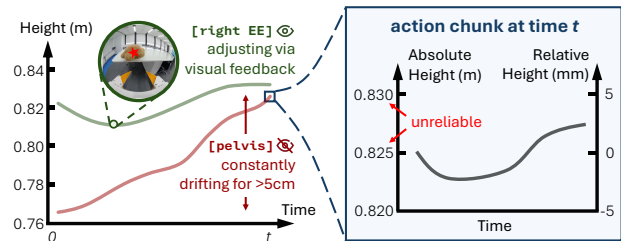


Figure 5. **Mitigating drift in non-vision-grounded keypoints.** **Left:** Trajectories during a doll-grasping task. The “sighted” gripper (green) remains anchored via visual feedback, whereas the “blind” pelvis (red) suffers from open-loop drift (> 5 cm) over time. **Right:** Decomposition of the action 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).

289 **Relative pose tracking for non-vision-grounded key-**
 290 **points.** Keypoint configuration is a critical design space for
 291 the policy interface. While previous frameworks often rely
 292 solely on gripper poses in the world frame [6, 12, 13], this
 293 is insufficient for whole-body manipulation. Theoretically,
 294 the HuMI demonstration data provides poses for full body
 295 keypoints; however, a trade-off exists between observabil-
 296 ity and the number of controlled keypoints. Unlike grippers
 297 equipped with wrist-view cameras, keypoints such as the
 298 pelvis and feet are “blind”—they lack direct visual anchors.
 299 Consequently, they accumulate unrecoverable errors during
 300 inference. As illustrated in Fig. 5 (left), during a station-
 301 ary grasping task, the target pelvis height drifts significantly
 302 (> 5 cm), rendering the absolute transform an unreliable
 303 control signal. To mitigate this, we modify the tracking ob-
 304 jective for non-vision-grounded keypoints to track the rel-
 305 ative transform within the current action chunk rather than
 306 the absolute transform (Fig. 5 right). This approach enables
 307 flexible keypoint configurations; we use a 3-keypoint setup
 308 (grippers + base) by default and demonstrate that the sys-
 309 tem maintains robust performance even when scaled to 5

310 keypoints (adding feet).

311 3. Experiments

312 In this section, we empirically evaluate HuMI along three
313 key dimensions. Specifically, we aim to answer the follow-
314 ing questions:

- 315 1. **Whole-body manipulation capability.** Can HuMI learn
316 feasible whole-body skills from robot-free demonstra-
317 tions, achieve sufficient manipulation precision while re-
318 specting whole-body dynamics, and effectively coordi-
319 nate high-low level policy during fully autonomous exe-
320 cution?
- 321 2. **Generalization ability.** Do robot-free demonstrations
322 collected across varied environments enable the learned
323 policy to generalize to unseen environments and objects?
- 324 3. **Data-collection efficiency.** Can HuMI acquire whole-
325 body manipulation data efficiently and with a high ac-
326 ceptance rate, and does the collected dataset cover ver-
327 satile whole-body skills, including motions that are chal-
328 lenging to obtain for teleoperation?

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

335 4. Whole-Body Manipulation Capability

336 In the capability experiments, we focus on evaluating
337 HuMI’s whole-body manipulation capability. We de-
338 sign four representative tasks that target complementary
339 aspects—whole-body coordination, precise bimanual ma-
340 nipulation, high-speed dynamic motion, and long-range
341 loco-manipulation. All experiments use in-domain settings,
342 with tasks evaluated in the same environments and initial
343 robot–object configurations as data collection.

344 4.1. Learning Feasible Whole-Body Skills from 345 Robot-Free Demonstrations

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

352 **Task challenges.** This task poses challenges for hu-
353 manoid whole-body coordination. The robot must coordi-
354 nate nearly all joints to transition from an upright stance to
355 a single-knee kneel while keeping its center of mass within
356 a very narrow support polygon and maintaining balance. At
357 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 col-
lection, 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. Un-
der this setting, the success rate drops from $17/20 = 85\%$
rollouts to $1/10 = 10\%$, and the learned policy often pro-
poses kinematically inappropriate motions. As illustrated
in Fig. 6 (b), the robot frequently kneels with an exces-
sively 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
removing whole-body supervision. Following prior UMI-
style interfaces [6, 12, 13], the high-level policy communi-
cates only end-effector waypoints, and the low-level con-
troller is trained from the two sparse EE waypoints. The
success rate drops to $0/10$. The low-level controller strug-
gles to maintain a stable, coordinated whole-body motion
from these underspecified demonstrations and often con-
verges to kinematically or dynamically inappropriate solu-
tions.

4.2. Precise Humanoid Bimanual Manipulation

To evaluate HuMI’s capability for precise humanoid biman-
ual 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 biman-
ual 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 mo-
tion induce shear that degrades alignment, making this mo-
tion a stringent benchmark for precise humanoid bimanual
manipulation.

Performance. HuMI successfully completes the un-
sheathing in $17/20 = 85\%$ trials, with an average end-
effector tracking error of 15.7 mm (Fig. 6 (c)). The robot
secures precise grasps on the hilt and scabbard without slip-
page or collisions with the rack. It then maintains tight bi-
manual coordination so that the blade slides out smoothly
in a single continuous pull (Fig. 6 (a)).

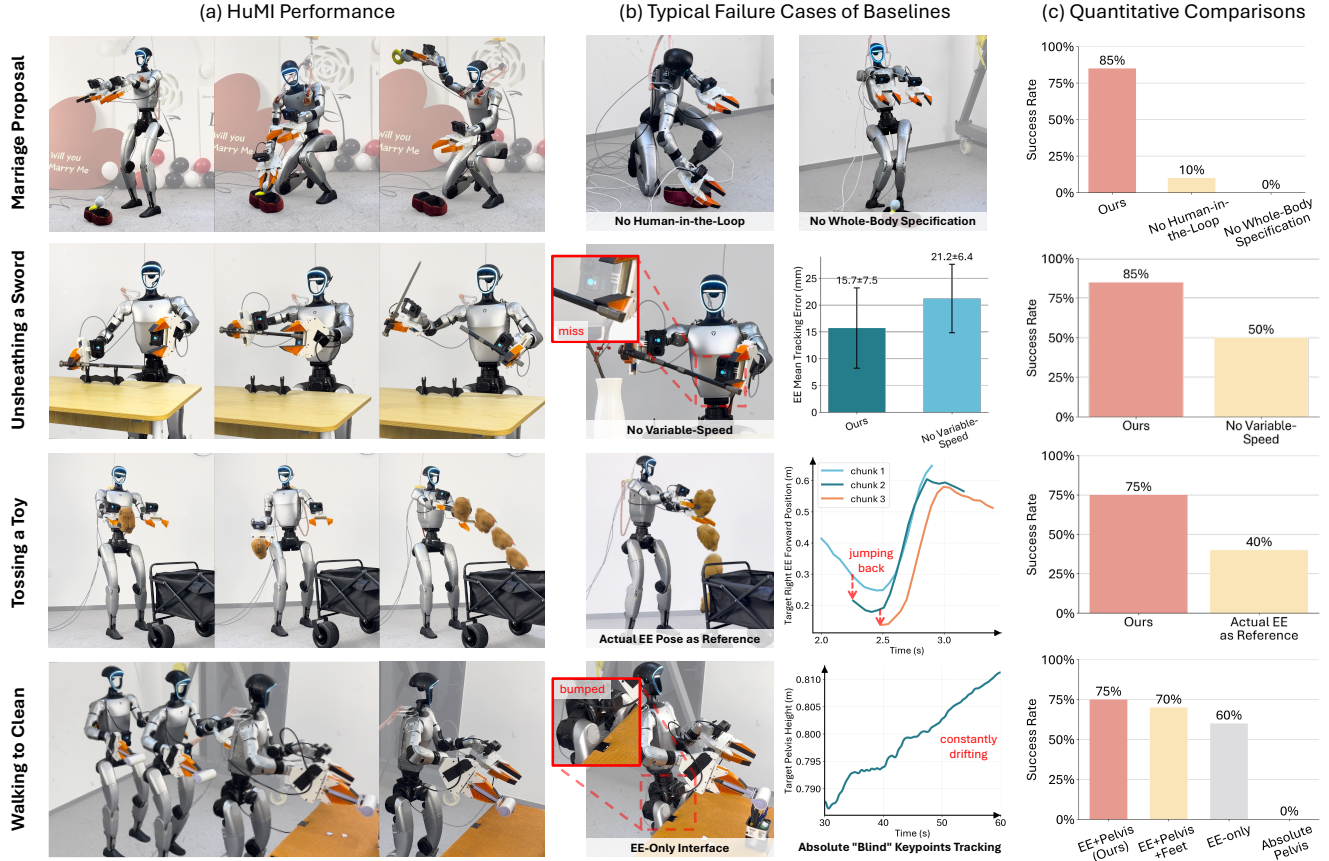


Figure 6. **Whole-body capability experiment results:** (a) **HuMI performance** across the four tasks. (b) **Typical failure modes** of the ablation baselines. The red dashed box highlights the failure behavior details. (c) **Success rate** for each task.

411 **Remove variable-speed augmentation.** We ablate the
 412 variable-speed augmentation by training the low-level policy
 413 on demonstrations replayed only at original human
 414 speeds. The success rate drops from 85% to $5/10 =$
 415 50% , with the average tracking error increasing to 21.2 mm
 416 (Fig. 6 (b)). The policy often fails to secure a reliable grasp
 417 or lacks the precise bimanual synchronization. These degrada-
 418 tions indicate that variable-speed augmentation is crucial
 419 for controller to resolve small tracking errors and improve
 420 bimanual coordination.

421 4.3. Temporally Coherent Dynamic Control

422 The third experiment investigates whether HuMI can ro-
 423 bustly and fluidly capture and transfer high-speed dynamic
 424 human motions to a humanoid robot. We consider **hu-
 425 manoid dynamic tossing** as the evaluation task. The robot
 426 must execute a temporally structured throwing motion, fea-
 427 turing a backward wind-up phase followed by a rapid forward
 428 swing that throws the object into a target container.

429 **Task challenges.** The robot must route momentum
 430 through many coupled joints, coordinating torso and arms
 431 for a backward wind-up and fast forward swing. These
 432 high-speed, continuous motions demand a tight control hi-

erarchy: the high-level policy outputs smooth, temporally
 coherent trajectory commands, while the low-level policy
 fluidly realizes them as coordinated whole-body motion.

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

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

4.4. 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 m away from the target desk, with its initial yaw offset sampled from $[-45^\circ, 45^\circ]$. 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.

5. 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 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.

6. 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

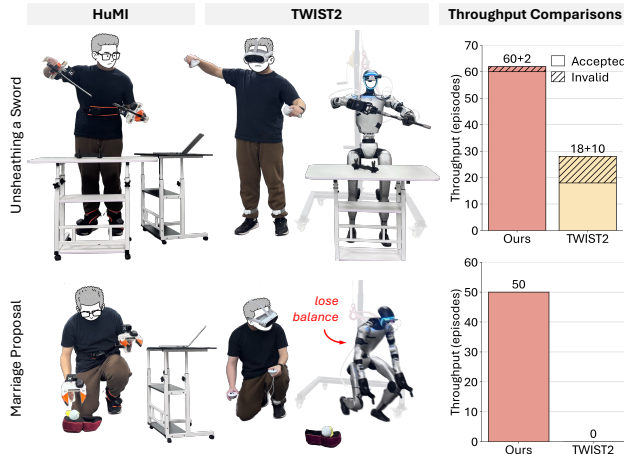


Figure 7. **Data collection throughput comparison.** **Left:** HuMI and TWIST2 workflows. **Right:** Number of episodes collected within 15 min; dashed segments denote invalid trajectories.

556 use the unsheathing task as a shared benchmark and run
 557 15 min collection sessions with both HuMI and TWIST2.
 558 For each system, we record the number of collected
 559 episodes and the acceptance rate. To ensure a fair compar-
 560 ison, all sessions are conducted by experienced users (more
 561 than 20 h with HuMI and more than 10 h with TWIST2).
 562 We define the acceptance rate as the fraction of episodes
 563 that are usable for downstream humanoid policy learning:
 564 an episode is acceptable only if the trajectory successfully
 565 completes the task and a policy trained on the full set of
 566 collected trajectories can replay this trajectory end-to-end.
 567 As summarized in Fig. 7 (upper), HuMI yields substan-
 568 tially higher throughput, collecting **62** episodes versus **28**
 569 with TWIST2, while also achieving a higher acceptance rate
 570 (**96.7%** vs. **64.3%**). HuMI’s streamlined, robot-free work-
 571 flow further reduces the average time per acceptable episode
 572 to **30.0%** of that of TWIST2, indicating that users can ob-
 573 tain dense, high-quality datasets much more quickly with
 574 our robot-free pipeline.

575 **Whole-body motion coverage.** To evaluate the ability
 576 to capture versatile behaviors, we again use the marriage
 577 proposal motion as a challenging target task. As shown in
 578 Fig. 7 (bottom), HuMI successfully collected **50** demon-
 579 strations within 15 minutes with a **100%** acceptance rate, aver-
 580 aging just **18 s** per episode. Conversely, TWIST2 failed to
 581 produce any usable demonstrations, as the teleoperated hu-
 582 manoid cannot reliably realize the required deep kneeling
 583 and often lose stability. This underscores a key advantage
 584 of HuMI’s robot-free data collection: it can capture diverse,
 585 highly articulated whole-body motions that go beyond the
 586 control limitations of existing humanoid teleoperation set-
 587 ups, providing broad coverage of complex behaviors for
 588 downstream policy learning.

7. 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 shap-
 ing and domain randomization [30, 34, 40], while teleop-
 eration imposes challenges for managing balance and com-
 pensating tracking errors. Although recent human video ap-
 proaches [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 manipu-
 lation across diverse tasks with strong generalization to un-
 seen 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, lack-
 ing the capacity for complex whole-body coordination. We
 present the first robot-free data collection system specifi-
 cally for humanoid whole-body manipulation.

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

Despite its efficacy, limitations remain. First, the sys-
 tem 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.

629

References

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

[1] Joao Pedro Araujo, Yanjie Ze, Pei Xu, Jiajun Wu, and C. Karen Liu. Retargeting matters: General motion re-targeting for humanoid motion tracking. *arXiv preprint arXiv:2510.02252*, 2025. 2, 3

[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. 1, 2, 4, 8, 9

[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. 3

[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. 9

[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. 9

[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. 2, 3, 5, 6, 9, 13, 14, 17

[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. 2, 5, 17

[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. 3, 9

[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. 9

[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. 13

[11] GoPro, Inc. Hero10 black — waterproof action camera. <https://gopro.com/en/us/shop/cameras/hero10-black/CHDXH-101-master.html>, 2021. Accessed: 2026-01-29. 13, 14

[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. 2, 3, 5, 6, 9

[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. 2, 3, 5, 6, 9

[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. 9, 15

[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. 1, 9

[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. 2, 15

[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. 3

[18] HTC VIVE. Vive ultimate tracker – full-body tracking for standalone vr, 2024. Accessed: 2026-01-18. 3, 9, 13, 14

[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. 9

[20] Lucy Lai, Ann Zixiang Huang, and Samuel J Gershman. Action chunking as policy compression. *PsyArXiv*, 2022. 2

[21] Jialong Li, Xuxin Cheng, Tianshu Huang, Shiqi Yang, Rizhao Qiu, and Xiaolong Wang. Amo: Adaptive motion optimization for hyper-dexterous humanoid whole-body control. *arXiv preprint arXiv:2505.03738*, 2025. 1, 8, 9

[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. 1, 4, 9, 15

[23] Qiayuan Liao, Takara E Truong, Xiaoyu Huang, Yuman Gao, Guy Tevet, Koushil Sreenath, and C Karen Liu. Beyond-mimic: From motion tracking to versatile humanoid control via guided diffusion. *arXiv preprint arXiv:2508.08241*, 2025. 2, 4, 17

[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. 9

[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. 1, 9

[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. 2, 3

[27] Zhengyi Luo, Chen Tessler, Toru Lin, Ye Yuan, Tairan He, Wenli Xiao, Yunrong Guo, Gal Chechik, Kris Kitani, Linxi

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

- 744 Fan, et al. Emergent active perception and dexterity of sim-
745 ulated humanoids from visual reinforcement learning. *arXiv*
746 *preprint arXiv:2505.12278*, 2025. 9
- 747 [28] Zhengyi Luo, Ye Yuan, Tingwu Wang, Chenran Li, Sirui
748 Chen, Fernando Castañeda, Zi-Ang Cao, Jiefeng Li, David
749 Minor, Qingwei Ben, et al. Sonic: Supersizing motion track-
750 ing for natural humanoid whole-body control. *arXiv preprint*
751 *arXiv:2511.07820*, 2025. 2, 4, 17
- 752 [29] OptiTrack. Primex 41 motion capture camera, 2024. Ac-
753 cessed: 2026-01-18. 3
- 754 [30] Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba,
755 and Pieter Abbeel. Sim-to-real transfer of robotic control
756 with dynamics randomization. In *2018 IEEE international*
757 *conference on robotics and automation (ICRA)*, pages 3803–
758 3810. IEEE, 2018. 9
- 759 [31] PICO. Pico virtual reality — official website, 2023. Ac-
760 cessed: 2026-01-25. 3
- 761 [32] Ri-Zhao Qiu, Shiqi Yang, Xuxin Cheng, Chaitanya Chawla,
762 Jialong Li, Tairan He, Ge Yan, David J. Yoon, Ryan Hoque,
763 Lars Paulsen, Ge Yang, Jian Zhang, Sha Yi, Guanya Shi, and
764 Xiaolong Wang. Humanoid policy ~ human policy. *arXiv*
765 *preprint arXiv:2503.13441*, 2025. 9
- 766 [33] Francisco J Romero-Ramirez, Rafael Muñoz-Salinas, and
767 Rafael Medina-Carnicer. Speeded up detection of squared
768 fiducial markers. *Image and vision Computing*, 76:38–47,
769 2018. 13
- 770 [34] Fereshteh Sadeghi and Sergey Levine. Cad2rl: Real single-
771 image flight without a single real image. *arXiv preprint*
772 *arXiv:1611.04201*, 2016. 9
- 773 [35] Nur Muhammad Mahi Shafiullah, Anant Rai, Haritheja
774 Etukuru, Yiqian Liu, Ishan Misra, Soumith Chintala, and
775 Lerrel Pinto. On bringing robots home. *arXiv preprint*
776 *arXiv:2311.16098*, 2023. 9
- 777 [36] Shuran Song, Andy Zeng, Johnny Lee, and Thomas
778 Funkhouser. Grasping in the wild: Learning 6dof closed-
779 loop grasping from low-cost demonstrations. *IEEE Robotics*
780 *and Automation Letters*, 5(3):4978–4985, 2020. 9
- 781 [37] Zhi Su, Bike Zhang, Nima Rahmanian, Yuman Gao, Qiayuan
782 Liao, Caitlin Regan, Koushil Sreenath, and S Shankar Sastry.
783 Hitter: A humanoid table tennis robot via hierarchical plan-
784 ning and learning. *arXiv preprint arXiv:2508.21043*, 2025.
785 2
- 786 [38] Generalist AI Team. Gen-0: Embodied foundation mod-
787 els that scale with physical interaction. *Generalist AI Blog*,
788 2025. <https://generalistai.com/blog/preview-uqlxvb-bb.html>.
789 9
- 790 [39] RDT Team. Rdt2: Enabling zero-shot cross-embodiment
791 generalization by scaling up umi data, 2025. 9
- 792 [40] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Woj-
793 ciech Zaremba, and Pieter Abbeel. Domain randomization
794 for transferring deep neural networks from simulation to the
795 real world. In *2017 IEEE/RSJ international conference on*
796 *intelligent robots and systems (IROS)*, pages 23–30. IEEE,
797 2017. 9
- 798 [41] Unitree Robotics. Unitree G1: Humanoid Agent AI Avatar.
799 <https://www.unitree.com/g1>, 2024. Accessed:
800 2026-01-16. 3, 4, 9, 14
- [42] Valve Corporation. Steamvr. <https://store.steampowered.com/app/250820/SteamVR/>,
2016. Accessed: 2026-01-29. 13
- [43] Vicon Motion Systems. Valkyrie optical motion capture
cameras, 2024. Accessed: 2026-01-18. 3
- [44] Timo Von Marcard, Bodo Rosenhahn, Michael J Black, and
Gerard Pons-Moll. Sparse inertial poser: Automatic 3d hu-
man pose estimation from sparse imus. In *Computer graph-
ics forum*, pages 349–360. Wiley Online Library, 2017. 3
- [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. 3, 9
- [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. 1, 9
- [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. 2, 3
- [48] Ruihan Yang, Qinxu 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. 9
- [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: Imper-
ative, web-based 3d visualization in python. *arXiv preprint*
arXiv:2507.22885, 2025. 13
- [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. 9
- [51] Kevin Zakka. Mink: Python inverse kinematics based on
MuJoCo, 2025. 13
- [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. 1, 4, 8, 9, 15
- [53] Yanjie Ze, Zixuan Chen, Wenhao Wang, Tianyi Chen, Xi-
alin He, Ying Yuan, Xue Bin Peng, and Jiajun Wu. Gener-
alizable humanoid manipulation with 3d diffusion policies.
In *2025 IEEE/RSJ International Conference on Intelligent*
Robots and Systems (IROS), pages 2873–2880. IEEE, 2025.
849
- [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. 1, 2, 4, 6, 8, 9
- [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 hu-
manoid balance. *arXiv preprint arXiv:2505.07294*, 2025. 855
856
857
858

- 859 [56] Siheng Zhao, Yanjie Ze, Yue Wang, C Karen Liu, Pieter
860 Abbeel, Guanya Shi, and Rocky Duan. Resmimic:
861 From general motion tracking to humanoid whole-body
862 loco-manipulation via residual learning. *arXiv preprint*
863 *arXiv:2510.05070*, 2025. 2
- 864 [57] Tony Z Zhao, Vikash Kumar, Sergey Levine, and Chelsea
865 Finn. Learning fine-grained bimanual manipulation with
866 low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.
867 2, 17

868 Please visit the offline website to view the robot
869 rollout videos. To access the content, please open
870 appendix-website/index.html in a web browser.

871 Appendix Contents

| | | |
|-----|---|-----------|
| 872 | A Data Collection System Details | 13 |
| 873 | A.1. Hardware Components | 13 |
| 874 | A.2. IK Interface | 13 |
| 875 | A.3. Data Processing | 13 |
| 876 | A.4. Data Collection Protocol | 14 |

| | | |
|-----|------------------------------------|-----------|
| 877 | B Deployment Details | 14 |
| 878 | B.1. Gripper Design | 14 |
| 879 | B.2. Camera Setup | 14 |
| 880 | B.3. System Architecture | 14 |

| | | |
|-----|---------------------------------|-----------|
| 881 | C Additional Experiments | 14 |
|-----|---------------------------------|-----------|

| | | |
|-----|---|-----------|
| 882 | D Generalization Experiments Details | 15 |
|-----|---|-----------|

| | | |
|-----|--|-----------|
| 883 | E Low-Level Controller Training Details | 15 |
|-----|--|-----------|

| | | |
|-----|---|-----------|
| 884 | F High-Level Policy Training Details | 17 |
|-----|---|-----------|

885 A. Data Collection System Details

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



Figure 8. HuMI’s hardware setup.

890 A.1. Hardware Components

891 As illustrated in Fig. 8, our data collection system com-
892 prises the following components:

- 893 • **Two UMI [6] grippers:** We utilize UMI grippers
894 equipped with GoPro cameras [11] to record wrist-view
895 RGB observations and ArUco markers [10, 33] for grip-
896 per width tracking.

- 897 • **Five HTC VIVE Ultimate trackers [18]:** To capture 6-
898 DoF poses, we attach trackers to the two grippers, the
899 waist, and the feet. We modified the top cover of the orig-
900 inal UMI gripper design to accommodate a tracker mount.
901 Standard VIVE straps are used to secure the trackers to
902 the waist and feet.

A laptop running SteamVR [42] is required to interface with the trackers to record their poses. 903 904

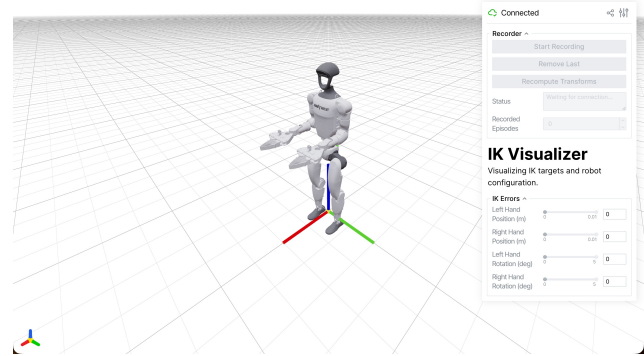


Figure 9. IK preview interface.

905 A.2. IK Interface

906 We developed a real-time IK preview interface to assist the
907 demonstrator in adapting their motions to be kinematically
908 feasible and task-compliant. Notably, to preserve strict spa-
909 tial relationships between the robot and the environment, we
910 use the original, unscaled motions (specifically, the track-
911 ers’ recorded $SE(3)$ transforms in the world frame) as IK
912 targets. Scaling is applied exclusively to the height of the
913 pelvis tracker.

914 The IK problem incorporates three subtasks: tracking the
915 translation and rotation of the five keypoints, avoiding self-
916 collisions, and maintaining a natural configuration via posture
917 regularization. We solve the IK problem in real-time using
918 Mink [51] and visualize the robot’s configuration using
919 Viser [49]. Fig. 9 depicts the user interface during data
920 collection.

921 A.3. Data Processing

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

- 925 1. **Synchronization:** We use the tracker timestamps as the
926 reference clock since the five trackers are natively syn-
927 chronized. To align the gripper videos with the tracker
928 data, we extract the gyroscope data embedded in the
929 video files by the GoPro cameras. We then align the
930 video and tracker timestamps by cross-correlating the
931 magnitudes of their angular velocities.

- 932 2. **Gripper width extraction:** Following Chi et al. [6], we
933 extract the gripper width from the recorded videos by
934 detecting the ArUco markers attached to the grippers.
935 3. **Data packaging:** The recorded data are packaged into
936 two subsets: (1) visual observations, gripper widths, and
937 keypoint trajectories for training the high-level policy;
938 and (2) keypoint trajectories paired with whole-body IK
939 solutions for training the low-level controller.

940 A.4. Data Collection Protocol

941 The following is the step-by-step protocol for collecting
942 demonstrations in a new scene:

- 943 1. **Mapping setup:** For each new scene, the demonstrator
944 must follow VIVE’s prompts to build a tracking map of
945 the environment; this typically takes 1–2 minutes.
946 2. **Calibration and synchronization:** The demonstrator
947 may optionally perform gripper width calibration and
948 GoPro timestamp synchronization following Chi et al.
949 [6].
950 3. **Demonstration:** The demonstrator repeatedly performs
951 the task within the scene. We use the start and stop
952 times of the recorded videos to delineate episodes. The
953 demonstrator uses GoPro’s voice commands to control
954 video recording, while the tracker data recording is controlled
955 via the IK interface GUI (Fig. 9). Thanks to the
956 synchronization step described in Sec. A.3, strict alignment
957 of start/stop times between the videos and tracker
958 recordings is not required. During the demonstration, the
959 demonstrator adapts their motions based on the real-time
960 IK preview.

961 B. Deployment Details

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

970 B.1. Gripper Design

971 To endow the Unitree G1 with manipulation capabilities
972 while minimizing distal mass, we developed a custom
973 gripper adapted from the UMI hardware interface [6].
974 We replaced the original spring-trigger mechanism with a
975 direct-drive transmission inspired by the Wild LMA design.
976 Specifically, we utilized the robot’s existing wrist yaw motor
977 to actuate the gripper, engineering a custom master gear
978 that mates precisely with the motor’s spline. This design
979 allows us to remove the stock rubber hands and clamping
980 mechanisms, securing the new mount via the original screw
981 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 en-
sure compatibility with our training data.

985 B.2. Camera Setup

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

993 B.3. System Architecture

994 As outlined in the main text, our control framework com-
995 prises a hierarchical structure: a high-level manipulation
996 policy and a low-level whole-body controller. Communi-
997 cation between the external workstation and the robot’s on-
998 board computer is established via a local wireless network
999 using ZeroMQ.

High-Level Policy The high-level policy runs on the
workstation at 5 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.

Low-Level Controller The low-level whole-body con-
troller executes directly on the robot’s onboard computer
at 50 Hz. Upon receiving the target keypoints and gripper
commands, it computes the necessary joint position com-
mands, which are executed by the robot’s built-in PD con-
troller. 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 di-
rectly from onboard sensors and streamed back to the high-
level policy.

1018 C. Additional Experiments

1019 In this section, we investigate the necessity of our adaptive
1020 end-effector (EE) reward. We aim to address a critical ques-
1021 tion raised in the methodology: whether naively prioritizing
1022 EE tracking precision in the current motion tracking frame-
1023 work would compromise whole-body stability. To validate
1024 this, we employ the **squat and pick up a bottle** task as
1025 a benchmark. This task is representative as it demands a
1026 seamless synergy between high-precision manipulation and
1027 whole-body dynamic balance: the humanoid must tightly
1028 coordinate its lower body and torso to maintain stability

1029 during the deep squat, while simultaneously achieving suf-
1030 ficient EE precision to execute a successful near-ground
1031 grasp and pick-up.

1032 **Remove adaptive end-effector reward.** To instantiate
1033 the baseline, we establish a baseline that enforces a **naive**
1034 **tight tracking constraint** throughout the entire motion.
1035 Specifically, we replace the adaptive end-effector (EE) re-
1036 ward with a fixed formulation:

$$1037 \quad r_{EE} = \sum_{\chi \in \{\mathbf{p}, \theta\}} r(\bar{e}_{\chi}^{EE}, \sigma_{\chi}^{\text{fixed}}),$$

1038 where we set $\sigma_{\chi}^{\text{fixed}}$ to the minimum tolerance used in our
1039 adaptive reward ($\sigma_{\mathbf{p}}^{\text{fixed}} = 0.01$ m, $\sigma_{\theta}^{\text{fixed}} = 5^{\circ}$), and fix the
1040 weight to $\omega_{EE}^{\text{fixed}} = \omega_{EE}^{\text{max}} = 0.5$. This configuration en-
1041 forces a uniformly tight EE tracking constraint throughout
1042 the motion, regardless of motion phase or modality. Under
1043 this setting, the success rate drops from $17/20 = 85\%$ to
1044 $5/10 = 50\%$. Typical failures occur during the deep-squat
1045 phase: the humanoid either loses balance and falls, or strug-
1046 gles to settle into an unnatural, marginally stable pose that
1047 prevents it from subsequently manipulation. This degrada-
1048 tion highlights the critical role of the adaptive EE reward.
1049 By dynamically shaping the reward landscape according to
1050 motion modalities, our method balances the need of whole-
1051 body stability and manipulation precision throughout train-
1052 ing. In contrast, permanently prioritizing EE precision
1053 forces the policy to neglect essential whole-body dynamics,
1054 compromising stability during whole-body dynamic phases
1055 like squatting.

1056 D. Generalization Experiments Details

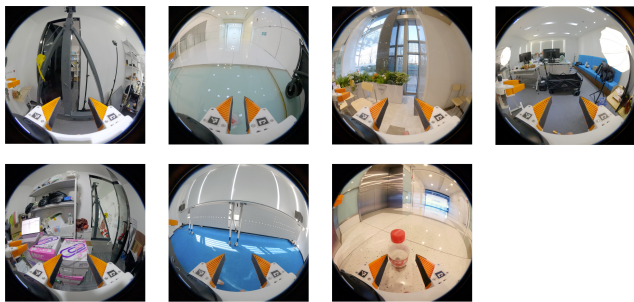


Figure 10. Training environments.

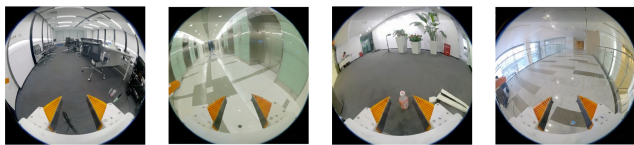


Figure 11. Testing environments.

1057 **Environment details.** Fig. 10 visualizes the seven train-
1058 ing environments used for our policy in Sec. 5. We collected

50 demonstrations in each environment, resulting in a to-
tal 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 ob-
jects used for our policy in Sec. 5. Fig. 13 visualizes the 6
testing objects.



Figure 12. Training objects.



Figure 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 observa-
tion is $O_t^{\text{tea}} = [s_t^{\text{tea}}, a_{t-1}^{\text{tea}}, c_t^{\text{tea}}]$, 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 ve-
locity $\mathbf{q}_t, \dot{\mathbf{q}}_t$, the base angular velocity ω_t , and the base
gravity vector g_t projected into the body frame; a_{t-1}^{tea} is
the previous action. The whole-body command $c_t^{\text{tea}} =$
 $[\mathbf{q}_t^{\text{ref}}, \dot{\mathbf{q}}_t^{\text{ref}}, \mathbf{p}_t^{\text{ref}}, \theta_t^{\text{ref}}, \mathbf{p}_t^{\text{ref}} - \mathbf{p}_t, \theta_t^{\text{ref}} \ominus \theta_t]$ contains the refer-
ence 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, θ_t .

For the student policy, the observation is $O_t^{\text{stu}} =$
 $[s_{t-25:t}^{\text{stu}}, a_{t-26:t-1}^{\text{stu}}, c_t^{\text{stu}}]$, where $s_t^{\text{stu}} = [\mathbf{q}_t, v_t, \omega_t, g_t]$ de-
notes the same real-world state features as above, and we
include a history of 25 steps of states and past actions to pro-
vide temporal context. The student command strictly aligns

Table 1. Rewards used in low-level controller training.

| Reward Term | Equation | Weight |
|---|---|-----------------------|
| <i>Tracking Rewards</i> | | |
| 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 \rightarrow 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 \rightarrow 0.5 |
| <i>Regularization Penalty</i> | | |
| Action rate | $\ a_t - a_{t-1}\ _2^2$ | -5×10^{-2} |
| Joint limits | $\sum_j \mathbb{I}(q_j \notin (q_j^{\text{min}}, q_j^{\text{max}}))$ | -10.0 |
| Undesired contact | $\sum_{i \notin \{\text{ankles, knees, hips}\}} \mathbb{I}(\ \mathbf{F}_i\ _2 > 1.0 \text{ N})$ | -0.1 |

Table 2. Domain randomization used in low-level controller training.

| Term | Description | Sampling Range |
|-------------------------------|---|--|
| <i>Physical randomization</i> | | |
| 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 | $\mu_d \sim \mathcal{U}[0.3, 1.2]$ |
| Resitution | randomize resitution 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 | $\Delta x \sim \mathcal{U}[\pm 0.025], \Delta y, z \sim \mathcal{U}[\pm 0.05]$ |
| End-effector mass | randomize the mass of end-effector | $m \sim \mathcal{U}[0.75, 1.25] \cdot m$ |
| <i>Velocity perturbations</i> | | |
| Push robot | periodically push the robot every Δt time | $\Delta t \sim \mathcal{U}[4, 6], v_{x,y} \sim \mathcal{U}[\pm 0.5], \omega_{\text{yaw}} \sim \mathcal{U}[\pm 0.78]$ |
| <i>Reset perturbations</i> | | |
| Base position | perturb base position at the reset time | $\Delta p_{x,y} \sim \mathcal{U}[\pm 0.05], \Delta p_z \sim \mathcal{U}[0.0, 0.05]$ |
| Base orientation | perturb base orientation at the reset time | $\Delta \omega_{p,r} \sim \mathcal{U}[\pm 0.1], \Delta \omega_{\text{yaw}} \sim \mathcal{U}[\pm 0.2]$ |
| Base linear velocity | perturb base linear velocity at the reset time | $\Delta v_{x,y} \sim \mathcal{U}[\pm 0.05], \Delta v_z \sim \mathcal{U}[\pm 0.2]$ |
| Base rotation velocity | perturb base rotation velocity at the reset time | $\Delta \omega_{r,p} \sim \mathcal{U}[\pm 0.52], \Delta \omega_{\text{yaw}} \sim \mathcal{U}[\pm 0.78]$ |
| 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_{reset} | $t_{\text{reset}} \sim t_{\text{reset}} + \mathcal{U}[-0.05, 0.05]$ |
| <i>Speed randomization</i> | | |
| Variable-speed augmentation | sample a play speed every 0.02s | $s \sim \mathcal{N}_{[0.25, 1.25]}(\mu = 1.0, \sigma^2)$ |

1089 with the high-level policy output, $c_t^{\text{stu}} = [c_{\mathcal{T}_t}^{\text{EE}}, c_{\mathcal{T}_t}^{\text{blind}}]$, where
 1090 each component contains a list of 10 waypoints sampled
 1091 over the next 2 s:

$$1092 \mathcal{T}_t = \{t_k = t + k \cdot \Delta t \mid k = 1, \dots, 10\}, \quad \Delta t = \left\lfloor \frac{2}{10 \delta t} \right\rfloor,$$

1093 and the control time step is $\delta t = 1/50$ s.

1094 The end-effector command at time step t_k is $c_{t_k}^{\text{EE}} =$

$[\mathbf{p}_{t_k}^{\text{ref}}, \theta_{t_k}^{\text{ref}}, \mathbf{p}_{t_k}^{\text{ref}} - \mathbf{p}_{t_k}, \theta_{t_k}^{\text{ref}} \ominus \theta_{t_k}]$, which includes the lo- 1095
 calized reference end-effector position and orientation 1096
 $\mathbf{p}_{t_k}^{\text{ref}}, \theta_{t_k}^{\text{ref}}$, together with the position and orientation deltas, 1097
 $\mathbf{p}_{t_k}^{\text{ref}} - \mathbf{p}_{t_k}$ and $\theta_{t_k}^{\text{ref}} \ominus \theta_{t_k}$. For “blind” points such 1098
 as the pelvis or feet, we instead use relative displacements 1099
 with respect to the current reference pose: $c_{t_k}^{\text{blind}} =$ 1100
 $[\mathbf{p}_{t_k}^{\text{ref}} - \mathbf{p}_t^{\text{ref}}, \theta_{t_k}^{\text{ref}} \ominus \theta_t^{\text{ref}}]$, which keeps the command *relative* 1101

1102 *within each action chunk*, independent of the absolute posi-
1103 tion.

1104 **Reward design.** Following prior works [23, 28], we de-
1105 compose the low-level reward into a tracking term r_{tracking}
1106 and a regularization term r_{penalty} . Table 1 summarizes all
1107 reward terms used for training the low-level controller. For
1108 the adaptive end-effector tracking rewards, we compute the
1109 tolerance parameter σ_χ by linearly interpolating between
1110 $[\sigma_\chi^{\min}, \sigma_\chi^{\max}]$ as a function of the commanded end-effector
1111 velocity $v_{\text{EE}}^{\text{ref}}$:

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

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

1113 For the adaptive end-effector position tracking rewards,
1114 we set $\sigma_{\text{p}}^{\min} = 0.01$ m and $\sigma_{\text{p}}^{\max} = 0.1$ m. For the adaptive
1115 end-effector rotation tracking rewards, we set $\sigma_{\theta}^{\min} = 5^\circ$
1116 and $\sigma_{\theta}^{\max} = 20^\circ$. For both rewards, we set $v^{\max} = 0.1$ m/s
1117 and $v^{\min} = 0.05$ m/s. Their reward weights are linearly
1118 increased from 0.0 to 0.5, and σ_{p}^{\min} is decreased from 0.1 m
1119 to 0.01 m over training steps 10,000 to 15,000.

1120 **Domain randomization.** Table 2 summarizes the do-
1121 main randomizations used for low-level controller training,
1122 grouped into physical, velocity, reset, and speed randomiza-
1123 tion.

1124 The *reset pose shift* targets the mismatch between the
1125 command trajectory and the robot state at deployment. Dur-
1126 ing training, the low-level policy tracks a replayed com-
1127 mand that ignores the current state, while at test time the
1128 high-level policy issues online commands. When facing
1129 the temporal misalignment or sudden corrections, low-level
1130 controller may fail due to the out-of-distribution issue. To
1131 expose the policy to such cases, in each episode we sample
1132 a reset time t_{reset} along the demonstration and initialize the
1133 robot to the pose at $t_{\text{reset}} + \mathcal{U}[-0.05, 0.05]$. This makes the
1134 robot start off from the reference yet still follow the same
1135 command trajectory.

1136 For *variable-speed augmentation*, every 0.02 s we sam-
1137 ple a speed scaling factor $s \sim \mathcal{N}_{[0.25, 1.25]}(1.0, \sigma^2)$ and use
1138 it to scale the reference time. The standard deviation σ is
1139 linearly increased from 10^{-4} to 1.0 between training steps
1140 10,000 and 15,000, after which it is kept at $\sigma = 1.0$.

1141 F. High-Level Policy Training Details

1142 We employ Diffusion Policy [7] as the backbone for high-
1143 level manipulation. Notably, this component is designed
1144 with modularity in mind; alternative architectures, such as
1145 ACT [57], can be seamlessly substituted.

1146 **Hyperparameters** The training hyperparameters remain
1147 consistent across all experimental tasks and are summarized
1148 in Table 3.

Table 3. **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× |
| 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 |

1149 **Observation and Action Spaces** The high-level poli-
1150 cy conditions on both visual and proprioceptive observa-
1151 tions. Visual inputs comprise RGB images captured by
1152 two gripper-mounted cameras, each with a resolution of
1153 224×224 pixels. Proprioceptive observations consist of
1154 the robot’s lower-body joint angles. The action space de-
1155 fines the desired positions and rotations of the keypoints,
1156 along with gripper widths. We adopt the same action pa-
1157 rameterization as UMI [6].

1158 **Training Data** We utilize 100 demonstrations collected
1159 in a single environment for the capability tasks. These tasks
1160 include marriage proposal, unsheathing a sword, tossing a
1161 toy, and walking to clean a table. For the task of squatting to
1162 pick up a bottle, we use 350 demonstrations collected across
1163 seven distinct environments (7×50).