



# CLONE: Closed-Loop Whole-Body Humanoid Teleoperation for Long-Horizon Tasks

Yixuan Li<sup>\*,1,2</sup>

lyx@bit.edu.cn

Tengyu Liu<sup>2,6</sup>

liutengyu@bigai.ai

Yutang Lin<sup>\*,3,4</sup>

yutang.lin@stu.pku.edu.cn

Jieming Cui<sup>2,3</sup>

cuijieming@stu.pku.edu.cn

Wei Liang<sup>1</sup>

liangwei@bit.edu.cn

Xixin Zhu<sup>3,5</sup>

yixin.zhu@pku.edu.cn

Siyuan Huang<sup>2,6</sup>

syhuang@bigai.ai

\* Equal contributors. <sup>1</sup> School of Computer Science and Technology, Beijing Institute of Technology

<sup>2</sup> State Key Laboratory of General Artificial Intelligence, Beijing Institute for General Artificial Intelligence (BIGAI)

<sup>3</sup> Institute for Artificial Intelligence, Peking University <sup>4</sup> Yuanpei College, Peking University

<sup>5</sup> Embodied Intelligence Lab, PKU-Wuhan Institute for Artificial Intelligence

<sup>6</sup> Joint Laboratory of Embodied AI and Humanoid Robots, BIGAI & UniTree Robotics

<https://humanoid-clone.github.io/>



Figure 1: **CLONE** employs an MoE-based policy with **closed-loop** error correction for humanoid teleoperation, enabling precise **whole-body coordination** and **long-horizon** task execution.

**Abstract:** humanoids offer unique potential for complex tasks requiring whole-body coordination, such as navigating environments while manipulating objects. However, current teleoperation systems face critical limitations: they decouple upper- and lower-body control to maintain stability, restricting natural coordination, and operate open-loop without real-time position feedback, leading to accumulated drift. The fundamental challenge is achieving precise, coordinated whole-body teleoperation over extended durations while maintaining accurate global positioning. Here we show that an MoE-based teleoperation system, **CLONE**, with closed-loop error correction enables unprecedented whole-body teleoperation fidelity, maintaining minimal positional drift over long-range trajectories using only head and hand tracking from an MR headset. Unlike previous methods that either sacrifice coordination for stability or suffer from unbounded drift, **CLONE** learns diverse motion skills while preventing tracking error accumulation through real-time feedback, enabling complex coordinated movements such as “picking up objects from the ground.” These results establish a new benchmark for whole-body humanoid teleoperation for long-horizon tasks.

**Keywords:** Humanoid; Whole-body teleoperation; Humanoid-scene interaction

## 1 Introduction

The ability to seamlessly coordinate whole-body movements while navigating complex environments represents one of humanity’s most remarkable capabilities [1, 2]. From squatting to retrieve objects from the ground to walking across rooms while carrying items, humans effortlessly integrate locomotion and manipulation in ways that remain challenging for robotic systems [3–7]. Humanoids (humanoids henceforth), with their human-like morphology, offer the promise of replicating these capabilities—potentially enabling applications from household assistance to operations in hazardous environments where human-like dexterity and mobility are essential [8–12].

However, realizing this potential requires solving a fundamental challenge: enabling intuitive, precise teleoperation that maintains coordination across the entire body over extended periods. Long-horizon tasks—such as navigating to distant locations while manipulating objects—demand not only moment-to-moment stability but also sustained accuracy in both movement execution and global positioning. Current teleoperation approaches fall short of these requirements, creating a significant capability gap between human operators and humanoids.

Recent advances in humanoid teleoperation and loco-manipulation [13–20] have made notable progress. Nevertheless, existing methods struggle with precise teleoperation over extended durations and fall short of the whole-body coordination necessary for humanoid-scene interaction. Two fundamental challenges persist in bridging this capability gap.

The first challenge centers on achieving **coordinated whole-body coordination**. Many systems decouple upper- and lower-body control for stability [18, 21], sacrificing the natural synergies required for fluid motion. While this separation provides safety, it fundamentally limits integrated actions such as reaching while walking or adjusting posture during manipulation. Alternative approaches that rely on motion capture data [13, 15, 19, 22–27] often emphasize stability at the cost of expressiveness, yielding conservative motions constrained by training data distributions. Moreover, these methods consistently overlook key factors like hand orientation that are critical for dexterous tasks, further restricting humanoids’ potential for sophisticated whole-body movements.

The second challenge involves **accumulated positional drift** over time due to the absence of real-time feedback about the robot’s actual position in the environment. Unlike wheeled robots with straightforward odometry, humanoids exhibit complex foot-ground interactions and non-holonomic dynamics that complicate accurate state estimation. Without closed-loop correction, small pose errors compound with each step, progressively degrading the operator’s spatial awareness and control authority, eventually leading to complete task failure. This drift problem becomes particularly acute during manipulation tasks that require precise positioning relative to environmental objects.

To tackle the above challenges in humanoid long-horizon tasks requiring whole-body coordination and accurate positioning, we present **CLONE**, a **closed-loop** whole-body teleoperation system combining learning-based coordination and real-time feedback correction. Our system employs a Mixture-of-Experts (MoE) architecture that learns to coordinate diverse motion skills while a LiDAR-based error correction mechanism prevents the accumulation of positional drift. Critically, **CLONE** requires only head and hand tracking from a single commercial Mixed Reality (MR) headset, making it practical for real-world deployment while achieving unprecedented fidelity in long-horizon tasks.

Our approach integrates three key components: (i) We develop an MoE framework that enables unified learning of diverse motion skills while maintaining natural upper- and lower-body coordination throughout execution. (ii) We implement closed-loop error correction using LiDAR odometry [28] and Apple Vision Pro (AVP) tracking to provide continuous global pose feedback and prevent drift accumulation. (iii) We curate a comprehensive dataset, **CLONED**, that augments AMASS [29] with hand orientations and additional motion-captured sequences, ensuring robust generalization to complex manipulation scenarios involving coordinated whole-body movements.

Our experiments demonstrate that **CLONE** enables capabilities previously unattainable with existing systems: whole-body coordination over long trajectories with minimal positional drift, complex co-

ordinated movements like object retrieval from ground level, and robust performance across diverse operator configurations and environmental conditions. Using only minimal input from a commercial MR headset, **CLONE** achieves improved tracking precision over existing open-loop approaches, opening new possibilities for practical humanoid applications in unstructured environments.

Our contributions are four-fold: (i) the first MoE-based framework for coordinated whole-body teleoperation that maintains natural movement synergies; (ii) a closed-loop system that solves the fundamental position drift problem in long-horizon tasks through real-time pose correction; (iii) a comprehensive dataset, **CLONED**, enabling robust learning of dexterous whole-body motions with proper hand orientation coverage; and (iv) extensive validation demonstrating substantial improvements in real-world humanoid-scene interaction capabilities across diverse scenarios.

## 2 Related Work

**Whole-Body Humanoid Teleoperation** Humanoid teleoperation enables robots to replicate human movements for complex tasks using motion capture systems [24, 17, 30], haptic devices [31–33], or virtual reality interfaces [15, 19, 34–36]. Control policies must balance robot stability with motion tracking fidelity. However, current methods struggle to reproduce the full diversity and fluidity of human motions [37], primarily due to monolithic MLP-based architectures that inadequately handle conflicting objectives in various types of motion (*e.g.*, walking *vs.* crouching) [38–40]. Although mixture-based models have shown promise in other domains [41–44], their application to humanoid teleoperation remains underexplored. In this work, we leverage an MoE framework for adaptive learning and to represent various motion patterns within a unified policy.

**Long-Horizon Loco-Manipulation** Long-horizon task execution [45] has been extensively studied for fixed-base arms [46–49], mobile manipulators [50–54], and aerial manipulators [55], typically in structured settings. In contrast, humanoid teleoperation remains limited to short-horizon motion replication [19, 23, 24], operating open-loop due to difficulties in real-time global state estimation for bipedal systems. Although recent advances in odometry have improved state tracking for legged robots [56, 57], their application to long-horizon humanoid control remains largely unexplored. To bridge this gap, we integrate LiDAR odometry into our teleoperation framework to enable closed-loop error correction and significantly reduce accumulated drift.

**Datasets for Training Humanoids** Large-scale motion capture (MoCap) datasets [29, 58] have been instrumental in training humanoid control policies [15, 19, 23, 59]. Even after augmenting the datasets with generative models [24], these datasets were still confined primarily to animation and graphics [60, 25] rather than robotics applications. While they contain semantically distinct actions (*e.g.*, waving, hugging, drinking), they underrepresent the kinematic configurations and dynamic transitions required for robust, generalizable controller training in real-world scenarios. To address these limitations, we introduce **CLONED** by augmenting AMASS [29] through motion editing and collecting additional human MoCap data, specifically tailored for humanoid controllers. This expansion increases coverage of motions and transitions relevant to humanoid control tasks.

## 3 The CLONE Framework

Our teleoperation framework, **CLONE**, captures a minimal set of control signals from the teleoperator, consisting solely of the 6D poses (position and orientation) of the head and both wrists, tracked using an AVP headset. These three points (see also Fig. 2) serve as the complete control interface, providing an intuitive yet powerful means of directing the humanoid’s whole-body motion while maintaining a simple setup that requires no additional hardware or complex calibration procedures.

**CLONE** addresses two fundamental challenges through complementary components. First, we develop a teacher-student policy learning approach that transforms these sparse control signals into coordinated whole-body movements (Sec. 3.1). Second, we implement a closed-loop error correction mechanism that maintains positional accuracy during extended operation (Sec. 3.2). The system

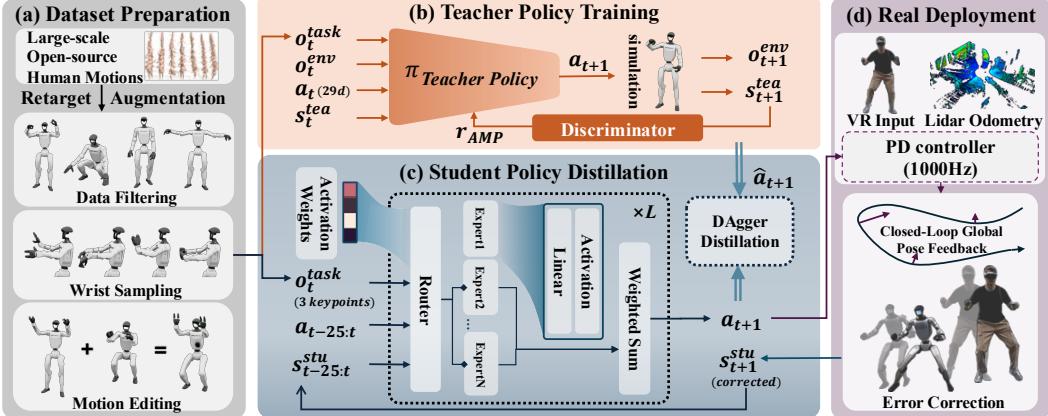


Figure 2: **The CLONE framework.** (a) CLONED curates and augments retargeted AMASS [29] data through motion editing to introduce diverse humanoid motions and detailed hand movements. (b) A teacher policy is trained using privileged information, including full robot state and environmental context. (c) An MoE network serves as the student policy, distilled from the teacher to operate with real-world observations only. (d) For real-world deployment, we integrate LiDAR odometry to obtain real-time humanoid states, enabling closed-loop error correction during teleoperation.

is supported by carefully designed reward structures and randomization techniques (Sec. 3.3) and trained on a newly curated dataset, **CLONED**, that ensures robust generalization (Sec. 3.4). Additional details of the implementation are provided in Appendix C.

### 3.1 Policy Learning

We employ a teacher-student training strategy for the teleoperation policy, following the overall framework of OmniH2O [19] (see Appendix A.1 for the problem formulation). This approach first trains a teacher policy with privileged information, then distills this knowledge into a student policy that operates using only real-world observations.

**Teacher Policy Training** The teacher policy  $\pi_{\text{tea}}$  is implemented as an Multi-Layer Perceptrons (MLP) that leverages comprehensive state information unavailable on real robots. At each timestep  $t$ , it processes observations  $\mathbf{o}_t^{\text{tea}} = [\mathbf{s}_t^{\text{tea}}, \mathbf{o}_t^{\text{task}}, \mathbf{a}_t, \mathbf{o}_t^{\text{env}}]$  and outputs target joint positions  $\mathbf{a}_{t+1} \in \mathbb{R}^{29}$  for PD control. The privileged states  $\mathbf{s}_t^{\text{tea}} = [\mathbf{p}_t, \theta_t, \mathbf{v}_t, \omega_t]$  include joint angular positions  $\mathbf{p}_t$  and the 6D poses, linear velocities, and angular velocities  $\theta_t, \mathbf{v}_t, \omega_t$  of all robot links. Task observations  $\mathbf{o}_t^{\text{task}} = [\hat{\mathbf{p}}_{t+1} - \mathbf{p}_t, \hat{\theta}_{t+1} - \theta_t, \hat{\mathbf{v}}_{t+1} - \mathbf{v}_t, \hat{\omega}_{t+1} - \omega_t, \hat{\mathbf{p}}_{t+1}, \hat{\theta}_{t+1}]$  capture both reference motion (denoted by  $\hat{\cdot}$ ) and tracking errors between reference and current states. Environmental observations  $\mathbf{o}_t^{\text{env}}$  provide context including ground friction coefficient and robot mass distribution.

**Student Policy Distillation** The student policy must operate without privileged information, following  $a_{t+1} = \pi_{\text{stu}}(s_{t-25:t}^{\text{stu}}, a_{t-25:t}, o_t^{\text{task}})$ . The robot state sequence  $s_{t-25:t}^{\text{stu}}$  contains joint positions  $q$ , joint velocities  $\dot{q}$ , root angular velocity  $\omega^{\text{root}}$ , and root gravity vector  $g$  obtained from on-device IMU over the past 25 frames. Task observations  $o_t^{\text{task}}$  consist of  $\hat{p}_{t+1} - p_t, \hat{\theta}_{t+1}, \hat{h}_{t+1}, h_t$ , and  $\hat{h}_{t+1}$ , where  $p_t$  represents the 3D positions of head and two wrists obtained from LiDAR odometry and forward kinematics,  $\hat{p}_{t+1}$  and  $\hat{\theta}_{t+1}$  are target positions and velocities from reference motion, and  $h_t, \hat{h}_{t+1}$  represent current and target wrist orientations.

The key challenge lies in handling diverse motion patterns within a single policy. Walking requires different control strategies than crouching or reaching, yet traditional monolithic architectures struggle with these conflicting objectives. We address this through an MoE architecture, as shown in Fig. 2, which allows specialized processing for different motion types.

The MoE design consists of  $L$  layers, each comprising  $N$  experts that function as independent feed-forward sub-layers with distinct parameters. At each layer, a router dynamically selects which experts are activated based on the input, generating weight distributions over all experts. The layer

output combines the top- $k$  experts with highest routing weights:  $f = \sum_i^k w_i \cdot E_i(\cdot)$ , where  $w_i$  is the routing weight for the  $i$ -th selected expert and  $E_i(\cdot)$  is the output of the  $i$ -th expert. This design enables different experts to focus on distinct motion patterns. To prevent model collapse to only a few experts, we introduce a balancing loss that encourages uniform expert selection:

$$\mathcal{L}_{balance} = \sum_{l=1}^L \sum_{e=1}^N [\max(p_e - \frac{1+\epsilon}{N}, 0) + \min(\frac{1-\epsilon}{N} - p_e, 0)], \quad (1)$$

where  $p_e = \mathbb{E}[w_e]$  represents the expected activation probability of expert  $e$ , and  $\epsilon$  is a slack constant that allows slight deviations from perfect uniformity.

### 3.2 Closed-Loop Error Correction

Traditional humanoid teleoperation systems operate in an open-loop configuration, where small errors in position tracking accumulate over time, leading to significant drift during extended operations. This fundamental limitation becomes particularly problematic during long-horizon tasks that require sustained positional accuracy. To address this challenge, we implement a closed-loop error correction mechanism that continuously monitors and compensates for positional discrepancies between the teleoperator and the humanoid.

Our approach utilizes LiDAR odometry to maintain accurate global position estimates for both the humanoid and the teleoperator. We employ FAST-LIO2 [28], an algorithm that tightly couples IMU and LiDAR data through an iterated Kalman filter to provide robust real-time state estimation even during dynamic movements (see more details in [Appendix A.2](#)). This choice ensures reliable tracking performance across diverse motion patterns, from walking to complex manipulation tasks.

The system tracks global positions for both agents: the humanoid’s position  $p \in \mathbb{R}^3$  is computed from onboard sensors, while the teleoperator’s position  $\hat{p} \in \mathbb{R}^3$  is similarly tracked through a MR hardware equipped with a comparable odometry pipeline. The student teleoperation policy directly consumes the difference between  $p$  and  $\hat{p}$ , enabling it to generate actions that systematically reduce positional drift and maintain accurate correspondence between the operator and the humanoid.

### 3.3 Reward Design and Domain Randomization

We build upon the reward terms and domain randomizations from OmniH2O [19] as the foundation of our approach, with specific enhancements to address the challenges of real-world teleoperation. Detailed reward functions and domain randomization settings are provided in [Appendix B](#).

To enhance robustness against LiDAR odometry errors, we introduce a velocity-dependent Stochastic Differential Equation (SDE) noise model during training that reflects real-world error characteristics. For the head position  $\vec{p}_{\text{head}}$ , we define the randomized position  $\vec{P}_{\text{head}}$  as:

$$d\vec{P}_{\text{head}} = \dot{\vec{p}}_{\text{head}} dt + \left( \frac{\|\dot{\vec{p}}_{\text{head}}\|}{c_{\text{vel}}} + c_{\text{min}} \right) d\vec{W}, \quad (2)$$

where  $\vec{W}$  is a standard Wiener process, and  $c_{\text{vel}}$  and  $c_{\text{min}}$  are constants that scale the noise proportionally to movement speed and establish a minimum randomization level. This formulation mirrors real-world dynamics, where faster movements tend to produce greater odometry errors. We use forward kinematics to compute other body positions based on the randomized head position, while periodically resetting and constraining the maximum deviation to avoid unrealistic drift.

Since **CLONE** provides only upper-body references (head and wrists), we must generate appropriate lower-body behaviors without explicit guidance. To tackle this challenge, we employ an Adversarial Motion Priors (AMP) reward [61] to regularize lower-body movements and encourage natural, stable behavior. Through this combination of specialized domain randomization and reward design, **CLONE** learns to generate robust lower-body behaviors while maintaining precise upper-body control aligned with operator commands.

### 3.4 The CLONED Dataset

The training dataset **CLONED** comprises three complementary components to support robust whole-body teleoperation. These include: (i) an augmented AMASS [29] subset of 149 curated sequences featuring diverse pairings of upper- and lower-body movements, enhanced via targeted motion editing to increase compositional diversity and policy generalization; (ii) 14 custom sequences captured with an IMU-based Xsens MoCap system to fill coverage gaps, emphasizing continuous transitions and diverse upper-body poses critical for manipulation; and (iii) systematic hand orientation augmentation through procedurally generated 6D wrist targets, smoothed via Spherical Linear Interpolation (SLERP) to ensure coherent and natural hand motions for teleoperation.

## 4 Real-World Experiments

For the real-world experiments, we evaluated **CLONE** on a physical Unitree G1 humanoid through both qualitative demonstrations and quantitative measurements, demonstrating unprecedented whole-body motion fidelity and precise position tracking. Our experiments focused on two key capabilities: (i) the accuracy of whole-body motion tracking for diverse skills, and (ii) the precision of our closed-loop error correction mechanism for global position tracking during extended teleoperation. Together, these experiments validate both the technical performance of our approach and its practical applicability to real-world humanoid teleoperation.

**Global Position Tracking** To quantitatively assess our system’s ability to maintain accurate global positioning over extended distances, we designed a controlled path-following experiment. We established fixed initial positions for both the operator and robot, marked on the floor for reference. The operator then walked along a straight path toward a pre-determined target position 3m, 6m, and 8.9m away while teleoperating the robot. We measured the discrepancy between the robot’s final position and the expected position as the tracking error. Over ten experimental iterations, our system demonstrated a mean tracking error of 5.1cm, with a maximum deviation of 12.0cm recorded at the 8.9m distance, as illustrated in Figure 3. This exceptional positional accuracy validates that **CLONE**’s closed-loop error correction mechanism effectively mitigates drift and accumulated positional errors throughout extended teleoperation sessions. Furthermore, the experimental results across the 3m, 6m, and 8.9m distances demonstrate that our model maintains consistently robust tracking performance, with both mean and maximum errors remaining within acceptable tolerances across the 8.9m operational range. Statistical analysis of the humanoid’s tracking accuracy revealed no significant performance differences across the tested distances of 3m, 6m, and 8.9m (all  $p$ -values  $> 0.05$ ). Independent samples t-tests comparing distance groups yielded the following results: 3m versus 6m ( $t = 0.165, p = 0.871$ ), 6m versus 8.9m ( $t = 0.048, p = 0.963$ ). Levene’s tests confirmed homogeneity of variance across all comparisons ( $p > 0.05$ ). The mean tracking errors exhibited remarkable consistency across distances, ranging between 0.0325m and 0.0356m. These findings substantiate that the robot’s tracking performance remains statistically equivalent within the tested operational range, thereby demonstrating the efficacy of the implemented closed-loop correction mechanisms.

**Whole-Body Motion Tracking** As shown in Fig. 4, **CLONE** successfully enables real-time teleoperation across a diverse range of whole-body skills. The robot accurately tracks complex motions including arm waving, deep squatting, standing up from squatted positions, and even dynamic jumping. To our knowledge, this represents notable advancements in whole-body motion fidelity for

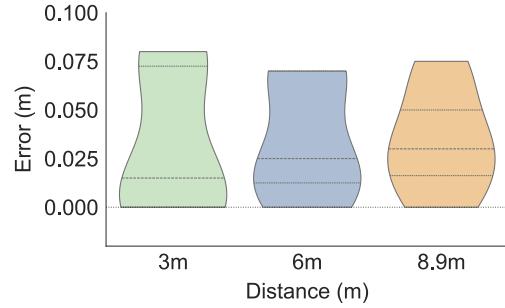


Figure 3: Long-Horizon Tracking Result of **CLONE** in Real-World.

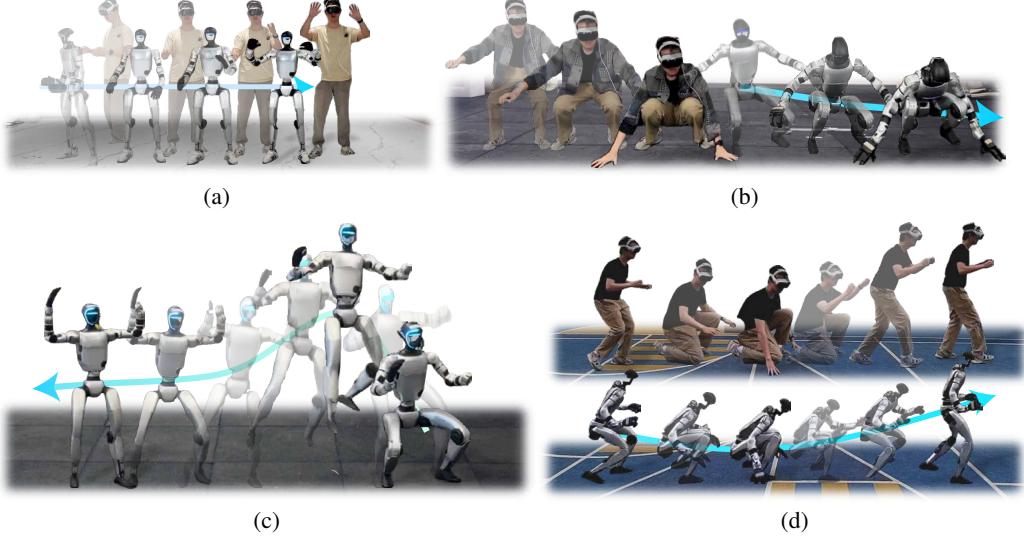


Figure 4: Qualitative results of motion tracking with **CLONE** on Unitree G1. Our method enables accurate tracking of diverse whole-body motions, including waving (a), squatting (b, d), and jumping (c).

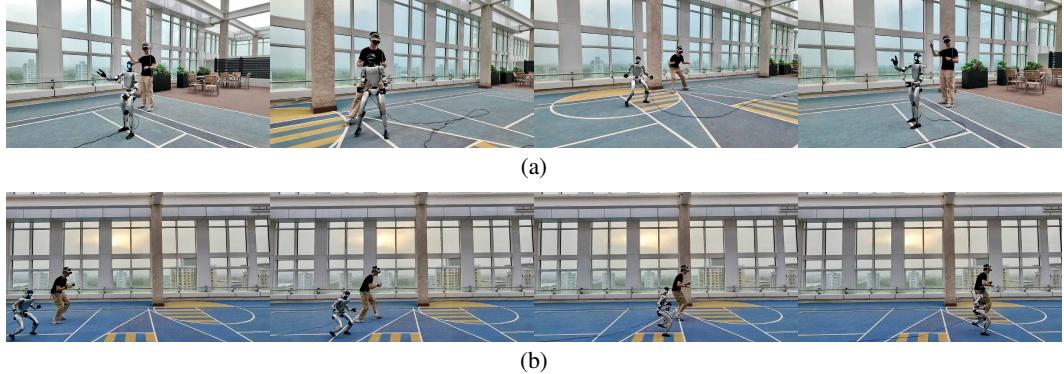


Figure 5: Qualitative results of long-horizon teleoperation. The humanoid accurately tracks both the human’s local pose and global translation over 15m and maintains its distance to its operator, as shown in (a, b).

real-time humanoid teleoperation, particularly for dynamic skills like jumping that require precise balance control and force application.

**Long-Horizon Mixed Navigation** To further validate our system’s performance in more complex scenarios, we conducted extended teleoperation sessions incorporating multiple movement types. As visualized in Fig. 5, we recorded a continuous teleoperation sequence where the operator traversed a complex path spanning over 15m, incorporating diverse locomotion patterns including forward walking, turning, side-stepping, and returning to the original position.

Throughout this extended sequence, the robot consistently tracked the operator’s movements with high accuracy and returned to its starting position with minimal drift. This demonstration highlights **CLONE**’s robustness for extended teleoperation sessions involving both locomotion and whole-body motion control—a capability essential for practical real-world applications of humanoids.

## 5 Simulations

In this section, we present a comprehensive evaluation of **CLONE** through extensive experiments conducted in simulation. Our evaluations span four key settings: reference motion tracking, diverse stance tracking, ablation studies, and expert activation analysis. These experiments are designed to: (i) quantitative benchmarking of motion tracking accuracy in the Isaac Gym simulation environ-

ment [62], (ii) assessment of robustness across diverse stance configurations, (iii) ablation studies examining key architectural decisions.

We evaluated **CLONE** on motion tracking tasks from **CLONED** using five metrics: success rate **SR** (%), mean per-keybody position error (MPKPE)  $E_{\text{mpkpe}}$  (mm), root-relative mean per-keybody position error (R-MPKPE)  $E_{\text{r-mpkpe}}$  (mm), average joint velocity error  $E_{\text{vel}}$  (mm/s), and hand orientation tracking error  $E_{\text{hand}}$ . Success rate (**SR**) represents the proportion of episodes where: (i) the robot maintains balance without falling, and (ii) the average per-keybody distance between the robot and reference motion remains below 1.5m across the three controlled joints. We defined the hand orientation tracking error as  $E_{\text{hand}} = 1 - \langle \hat{\mathbf{q}}, \mathbf{q} \rangle^2$ , where  $\hat{\mathbf{q}}$  and  $\mathbf{q}$  represent the reference and robot hand quaternions.

We compared **CLONE** with two ablated baselines, **CLONE**<sup>†</sup>, and **CLONE**<sup>\*</sup>. **CLONE**<sup>†</sup> employs an MLP as the student policy, resembling the OmniH2O baseline trained on our data and task. **CLONE**<sup>\*</sup> represents our **CLONE** model trained on OmniH2O data. Quantitative results in Tab. 1 reveal that both the MoE architecture and **CLONED** contribute significantly contribute significantly to accurate reference motion tracking. We provide qualitative results in Appendix D.1.

**Tracking Diverse Stances** To assess **CLONE**'s robustness across varying postures, we evaluated its performance in tracking motions with head heights from 1.2m (standing) to 0.6m (deep squatting) in 0.1m decrements. We generated these reference motions by systematically editing sequences from the **CLONED** dataset, creating unseen poses that challenge teleoperation systems.

As shown in Fig. 6, **CLONE** demonstrates an interesting trade-off: while underperforming baselines in absolute position accuracy (MPKPE), it consistently outperforms them in local metrics (R-MPKPE, velocity error, and hand orientation). This pattern suggests that **CLONE** prioritizes faithful reproduction of reference stances—particularly for challenging postures—sometimes at the expense of global positioning. All methods show increased tracking errors at lower heights, confirming the challenge of teleoperating robots in squatting postures.

## 6 Conclusion

We present a closed-loop teleoperation framework based on MoE architectures that addresses the fundamental challenge of tracking error accumulation in long-horizon humanoid control. The system leverages an augmented training dataset constructed through motion editing algorithms and MoCap data refinement processes.

Our approach integrates LiDAR-based odometry for enhanced spatial localization, implementing a stepwise error correction mechanism that maintains tracking accuracy throughout extended operation periods. The resulting system demonstrates versatility across a broad spectrum of tasks—encompassing complex manipulation behaviors, long-range navigation, and multi-phase task sequences—while simultaneously serving as a robust platform for collecting diverse demonstration data across various operational scenarios.

Table 1: **Motion tracking evaluation on the CLONED dataset.** We compare **CLONE** against ablation variants: **CLONE**<sup>†</sup> replaces MoE with an MLP as the student policy, and **CLONE**<sup>\*</sup> is trained on the original OmniH2O dataset instead of **CLONED**.

Method	<b>SR</b> $\uparrow$	$E_{\text{mpkpe}} \downarrow$	$E_{\text{r-mpkpe}} \downarrow$	$E_{\text{vel}} \downarrow$	$E_{\text{hand-rot}} \downarrow$
<b>CLONE</b> <sup>†</sup>	100%	113.97	35.55	245.11	4.73
<b>CLONE</b> <sup>*</sup>	100%	102.20	41.07	309.65	4.61
<b>CLONE</b>	100%	<b>87.84</b>	<b>33.30</b>	<b>227.17</b>	<b>3.61</b>

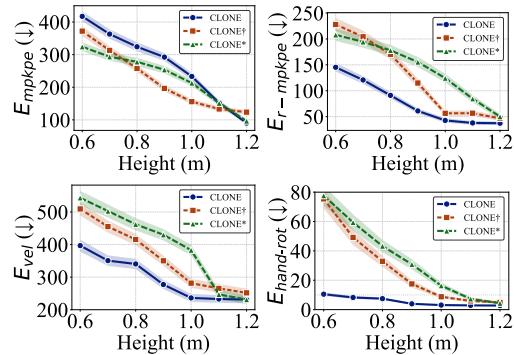


Figure 6: Motion tracking performance comparison between **CLONE** (blue solid), **CLONE**<sup>\*</sup> (green dashed), and **CLONE**<sup>†</sup> (red dashed). Lower values indicate better performance for all error metrics.

## 7 Limitation

While our approach demonstrates significant capabilities in humanoid teleoperation, several important limitations remain to be addressed in future work:

**Control Precision with Minimal Input** Our closed-loop error correction mechanism effectively reduces jitter and drift from accumulated errors during humanoid locomotion, enabling whole-body control with just three tracked points (head and hands) from a MR headset. However, this minimal input configuration inherently constrains the system’s stability in certain scenarios. Although our adversarial reward helps stabilize the lower body, fine-grained stability control remains challenging. Future work should explore additional sensing modalities or predictive algorithms to enhance stability while maintaining the simplicity of the interface.

**Dynamic Motion Capabilities** Despite our dataset augmentation strategies to increase motion diversity, the system exhibits reduced performance when executing highly dynamic movements such as jumping or single-leg hopping. These limitations stem from both the training data distribution and the inherent difficulty of balancing during such maneuvers. Addressing these constraints would require: (i) expanding the reference motion dataset to include more diverse dynamic sequences, (ii) developing specialized reward functions targeted at maintaining balance during rapid transitions, and (iii) potentially incorporating physics-based constraints into the policy to better model momentum dynamics.

## 8 Acknowledgments

We thank Unitree Robotics for providing essential hardware repairs.

This work is supported in part by the National Science and Technology Major Project (2022ZD0114900), the National Natural Science Foundation of China (62376031), the Beijing Nova Program, the State Key Lab of General AI at Peking University, the PKU-BingJi Joint Laboratory for Artificial Intelligence, and the National Comprehensive Experimental Base for Governance of Intelligent Society, Wuhan East Lake High-Tech Development Zone.

## References

- [1] B. Henze, M. A. Roa, and C. Ott. Passivity-based whole-body balancing for torque-controlled humanoid robots in multi-contact scenarios. *International Journal of Robotics Research (IJRR)*, 35(12):1522–1543, 2016.
- [2] P. M. Wensing and D. E. Orin. Improved computation of the humanoid centroidal dynamics and application for whole-body control. *International Journal of Humanoid Robotics*, 13(01):1550039, 2016.
- [3] L. Sentis and O. Khatib. Synthesis of whole-body behaviors through hierarchical control of behavioral primitives. *International Journal of Humanoid Robotics*, 2(04):505–518, 2005.
- [4] O. Khatib. A unified approach for motion and force control of robot manipulators: The operational space formulation. *IEEE Journal on Robotics and Automation*, 3(1):43–53, 2003.
- [5] T. Fukuda, P. Dario, and G.-Z. Yang. Humanoid robotics—history, current state of the art, and challenges. *Science Robotics*, 2(13):eaar4043, 2017.
- [6] A. Hereid, C. M. Hubicki, E. A. Cousineau, and A. D. Ames. Dynamic humanoid locomotion: A scalable formulation for hzd gait optimization. *IEEE Transactions on Robotics (T-RO)*, 34(2):370–387, 2018.

- [7] J. Cui, T. Liu, Z. Meng, J. Yu, R. Song, W. Zhang, Y. Zhu, and S. Huang. Grove: A generalized reward for learning open-vocabulary physical skill. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025.
- [8] Y. Tong, H. Liu, and Z. Zhang. Advancements in humanoid robots: A comprehensive review and future prospects. *IEEE/CAA Journal of Automatica Sinica*, 11(2):301–328, 2024.
- [9] Vecna Robotics. Battlefield Extraction-Assist Robot (BEAR). <https://robotsguide.com/robots/bear>, 2006.
- [10] U.S. Naval Research Laboratory. Autonomous Shipboard Humanoid (ASH). <https://www.navy.mil/Resources/Fact-Files/Display-FactFiles/Article/2160601/shipboard-autonomous-firefighting-robot-saffir/>, 2012.
- [11] AGIBOT Robotics. AGIBOT A2 Humanoid Robot. <https://www.agibot.com/products/A2>, 2024.
- [12] Z. Zhao, W. Li, Y. Li, T. Liu, B. Li, M. Wang, K. Du, H. Liu, Y. Zhu, Q. Wang, K. Althoefer, and S.-C. Zhu. Embedding high-resolution touch across robotic hands enables adaptive human-like grasping, 2024.
- [13] Z. Fu, Q. Zhao, Q. Wu, G. Wetzstein, and C. Finn. Humanplus: Humanoid shadowing and imitation from humans. In *Conference on Robot Learning (CoRL)*, 2024.
- [14] Y. Ze, Z. Chen, W. Wang, T. Chen, X. He, Y. Yuan, X. B. Peng, and J. Wu. Generalizable humanoid manipulation with improved 3d diffusion policies, 2024.
- [15] T. He, W. Xiao, T. Lin, Z. Luo, Z. Xu, Z. Jiang, C. Liu, G. Shi, X. Wang, L. Fan, and Y. Zhu. Hover: Versatile neural whole-body controller for humanoid robots, 2024.
- [16] J. Cui, T. Liu, N. Liu, Y. Yang, Y. Zhu, and S. Huang. Anyskill: Learning open-vocabulary physical skill for interactive agents. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- [17] S. Dafarra, K. Darvish, R. Grieco, G. Milani, U. Pattacini, L. Rapetti, G. Romualdi, M. Salvi, A. Scalzo, I. Sorrentino, et al. icub3 avatar system: Enabling remote fully immersive embodiment of humanoid robots. *Science Robotics*, 9(86):eadh3834, 2024.
- [18] Q. Ben, F. Jia, J. Zeng, J. Dong, D. Lin, and J. Pang. Homie: Humanoid loco-manipulation with isomorphic exoskeleton cockpit. In *Robotics: Science and Systems (RSS)*, 2025.
- [19] T. He, Z. Luo, X. He, W. Xiao, C. Zhang, W. Zhang, K. Kitani, C. Liu, and G. Shi. Omnih2o: Universal and dexterous human-to-humanoid whole-body teleoperation and learning. In *Conference on Robot Learning (CoRL)*, 2024.
- [20] T. Haarnoja, B. Moran, G. Lever, S. H. Huang, D. Tirumala, J. Humplik, M. Wulfmeier, S. Tunyasuvunakool, N. Y. Siegel, R. Hafner, et al. Learning agile soccer skills for a bipedal robot with deep reinforcement learning. *Science Robotics*, 9(89):eadi8022, 2024.
- [21] A. Matsiko. Humanoid robot learning of complex behaviors with llms. *Science Robotics*, 10(98):eadv4627, 2025.
- [22] T. He, Z. Luo, W. Xiao, C. Zhang, K. Kitani, C. Liu, and G. Shi. Learning human-to-humanoid real-time whole-body teleoperation. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2024.
- [23] X. Cheng, Y. Ji, J. Chen, R. Yang, G. Yang, and X. Wang. Expressive whole-body control for humanoid robots. In *Robotics: Science and Systems (RSS)*, 2024.

- [24] M. Ji, X. Peng, F. Liu, J. Li, G. Yang, X. Cheng, and X. Wang. Exbody2: Advanced expressive humanoid whole-body control, 2024.
- [25] N. Jiang, Z. Zhang, H. Li, X. Ma, Z. Wang, Y. Chen, T. Liu, Y. Zhu, and S. Huang. Scaling up dynamic human-scene interaction modeling. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- [26] N. Jiang, Z. He, H. Li, Y. Chen, S. Huang, and Y. Zhu. Autonomous character-scene interaction synthesis from text instruction. In *SIGGRAPH Asia Conference Papers*, 2024.
- [27] N. Jiang, H. Li, Z. Yuan, Z. He, Y. Chen, T. Liu, Y. Zhu, and S. Huang. Dynamic motion blending for versatile motion editing. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025.
- [28] W. Xu, Y. Cai, D. He, J. Lin, and F. Zhang. Fast-lio2: Fast direct lidar-inertial odometry. *IEEE Transactions on Robotics (T-RO)*, 38(4):2053–2073, 2022.
- [29] N. Mahmood, N. Ghorbani, N. F. Troje, G. Pons-Moll, and M. J. Black. AMASS: Archive of motion capture as surface shapes. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2019.
- [30] K. Darvish, Y. Tirupachuri, G. Romualdi, L. Rapetti, D. Ferigo, F. J. A. Chavez, and D. Pucci. Whole-body geometric retargeting for humanoid robots. In *International Conference on Humanoid Robots (Humanoids)*, 2019.
- [31] A. Brygo, I. Sarakoglou, N. Garcia-Hernandez, and N. Tsagarakis. Humanoid robot teleoperation with vibrotactile based balancing feedback. In *Haptics: Neuroscience, Devices, Modeling, and Applications*, 2014.
- [32] L. Peternel and J. Babič. Learning of compliant human–robot interaction using full-body haptic interface. *Advanced Robotics*, 27(13):1003–1012, 2013.
- [33] J. Ramos and S. Kim. Dynamic locomotion synchronization of bipedal robot and human operator via bilateral feedback teleoperation. *Science Robotics*, 4(35):eaav4282, 2019.
- [34] J. Chagas Vaz, D. Wallace, and P. Y. Oh. Humanoid loco-manipulation of pushed carts utilizing virtual reality teleoperation. In *International Mechanical Engineering Congress and Exposition*, 2021.
- [35] L. Penco, N. Scianca, V. Modugno, L. Lanari, G. Oriolo, and S. Ivaldi. A multimode teleoperation framework for humanoid loco-manipulation: An application for the icub robot. *IEEE Robotics & Automation Magazine*, 26(4):73–82, 2019.
- [36] S. Tachi, Y. Inoue, and F. Kato. Telesar vi: Telexistence surrogate anthropomorphic robot vi. *International Journal of Humanoid Robotics*, 17(05):2050019, 2020.
- [37] M. Moniruzzaman, A. Rassau, D. Chai, and S. M. S. Islam. Teleoperation methods and enhancement techniques for mobile robots: A comprehensive survey. *Robotics and Autonomous Systems*, 150:103973, 2022.
- [38] R. Huang, S. Zhu, Y. Du, and H. Zhao. Moe-loco: Mixture of experts for multitask locomotion, 2025.
- [39] S. Zhou, W. Zhang, J. Jiang, W. Zhong, J. GU, and W. Zhu. On the convergence of stochastic multi-objective gradient manipulation and beyond. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- [40] K. Darvish, L. Penco, J. Ramos, R. Cisneros, J. Pratt, E. Yoshida, S. Ivaldi, and D. Pucci. Teleoperation of humanoid robots: A survey. *IEEE Transactions on Robotics (T-RO)*, 39(3):1706–1727, 2023.

- [41] C. Yang, K. Yuan, Q. Zhu, W. Yu, and Z. Li. Multi-expert learning of adaptive legged locomotion. *Science Robotics*, 5(49):eabb2174, 2020.
- [42] Z. Xie, S. Starke, H. Y. Ling, and M. van de Panne. Learning soccer juggling skills with layer-wise mixture-of-experts. In *ACM SIGGRAPH / Eurographics Symposium on Computer Animation (SCA)*, 2022.
- [43] W. Song, H. Zhao, P. Ding, C. Cui, S. Lyu, Y. Fan, and D. Wang. Germ: A generalist robotic model with mixture-of-experts for quadruped robot. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2024.
- [44] G. Cheng, L. Dong, W. Cai, and C. Sun. Multi-task reinforcement learning with attention-based mixture of experts. *IEEE Robotics and Automation Letters (RA-L)*, 8(6):3812–3819, 2023.
- [45] C. R. Garrett, R. Chitnis, R. Holladay, B. Kim, T. Silver, L. P. Kaelbling, and T. Lozano-Pérez. Integrated task and motion planning. *Annual Review of Control, Robotics, and Autonomous Systems*, 4(1):265–293, 2021.
- [46] C. Wang, L. Fan, J. Sun, R. Zhang, L. Fei-Fei, D. Xu, Y. Zhu, and A. Anandkumar. Mimicplay: Long-horizon imitation learning by watching human play. In *Conference on Robot Learning (CoRL)*, 2023.
- [47] Z. Lin, Y. Chen, and Z. Liu. Hierarchical human-to-robot imitation learning for long-horizon tasks via cross-domain skill alignment. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2024.
- [48] H. Shi, H. Xu, S. Clarke, Y. Li, and J. Wu. Robocook: Long-horizon elasto-plastic object manipulation with diverse tools. *Conference on Robot Learning (CoRL)*, 2023.
- [49] Z. Zhao, Y. Li, W. Li, Z. Qi, L. Ruan, Y. Zhu, and K. Althoefer. Tac-man: Tactile-informed prior-free manipulation of articulated objects. *IEEE Transactions on Robotics (T-RO)*, 41: 538–557, 2024.
- [50] Z. Jiao, Z. Zeyu, W. Wang, D. Han, S.-C. Zhu, Y. Zhu, and H. Liu. Efficient task planning for mobile manipulation: a virtual kinematic chain perspective. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2021.
- [51] Z. Jiao, Z. Zeyu, X. Jiang, D. Han, S.-C. Zhu, Y. Zhu, and H. Liu. Consolidating kinematic models to promote coordinated mobile manipulations. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2021.
- [52] Z. Jiao, Y. Niu, Z. Zhang, S.-C. Zhu, Y. Zhu, and H. Liu. Planning sequential tasks on contact graph. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2022.
- [53] P. Zhi, Z. Zhang, M. Han, Z. Zhang, Z. Li, Z. Jiao, B. Jia, and S. Huang. Closed-loop open-vocabulary mobile manipulation with gpt-4v. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2024.
- [54] P. Zhi, P. Li, J. Yin, B. Jia, and S. Huang. Learning unified force and position control for legged loco-manipulation, 2025.
- [55] Y. Su, J. Li, Z. Jiao, M. Wang, C. Chu, H. Li, Y. Zhu, and H. Liu. Sequential manipulation planning for over-actuated unmanned aerial manipulators. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2023.
- [56] D. Wisth, M. Camurri, and M. Fallon. Vilens: Visual, inertial, lidar, and leg odometry for all-terrain legged robots. *IEEE Transactions on Robotics (T-RO)*, 39(1):309–326, 2022.

- [57] G. Ou, D. Li, and H. Li. Leg-kilo: Robust kinematic-inertial-lidar odometry for dynamic legged robots. *IEEE Robotics and Automation Letters (RA-L)*, 2024.
- [58] F. G. Harvey, M. Yurick, D. Nowrouzezahrai, and C. Pal. Robust motion in-betweening. *ACM Transactions on Graphics (TOG)*, 39(4):60–1, 2020.
- [59] L. Ma, Z. Meng, T. Liu, Y. Li, R. Song, W. Zhang, and S. Huang. Styleloco: Generative adversarial distillation for natural humanoid robot locomotion, 2025.
- [60] Z. Wang, Y. Chen, T. Liu, Y. Zhu, W. Liang, and S. Huang. Humanise: Language-conditioned human motion generation in 3d scenes. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- [61] X. B. Peng, Z. Ma, P. Abbeel, S. Levine, and A. Kanazawa. Amp: Adversarial motion priors for stylized physics-based character control. *ACM Transactions on Graphics (TOG)*, 40(4):1–20, 2021.
- [62] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin, A. Allshire, A. Handa, et al. Isaac gym: High performance gpu-based physics simulation for robot learning, 2021.