# BeyondMimic: From Motion Tracking to Versatile Humanoid Control via Guided Diffusion

Qiayuan Liao[1*†], Takara E. Truong[2†], Xiaoyu Huang[1†],

Yuman Gao[1], Guy Tevet[2], Koushil Sreenath[1‡], C. Karen Liu[2‡]

[1]University of California, Berkeley, CA 94720, USA.

[2]Stanford University, Stanford, CA 94305, USA.

[*]Corresponding author. Email: qiayuanl@berkeley.edu

[†]These authors contributed equally to this work; order decided by coin toss.

[‡]Equal advising; order mirrors the coin toss, with the other lab listed last.

**The human-like form of humanoid robots positions them uniquely to achieve the agility and versatility in motor skills that humans possess. Learning from human demonstrations offers a scalable approach to acquiring these capabilities. However, prior works either produce unnatural motions or rely on motion-specific tuning to achieve satisfactory naturalness. Furthermore, these methods are often motion- or goal-specific, lacking the versatility to compose diverse skills, especially when solving unseen tasks. We present BeyondMimic, a framework that scales to diverse motions and carries the versatility to compose them seamlessly in tackling unseen downstream tasks. At heart, a compact motion-tracking formulation enables mastering a wide range of radically agile behaviors, including aerial cartwheels, spin-kicks, flip-kicks, and sprinting, with a single setup and shared hyperparameters, all while achieving state-of-the-art human-like performance. Moving beyond the mere imitation of existing motions, we propose a unified la-**

1

**tent diffusion model that empowers versatile goal specification, seamless task switching, and dynamic composition of these agile behaviors. Leveraging classifier guidance, a diffusion-specific technique for test-time optimization toward novel objectives, our model extends its capability to solve downstream tasks never encountered during training, including motion inpainting, joystick teleoperation, and obstacle avoidance, and transfers these skills zero-shot to real hardware. This work opens new frontiers for humanoid robots by pushing the limits of scalable human-like motor skill acquisition from human motion and advancing seamless motion synthesis that achieves generalization and versatility beyond training setups.**

## Introduction

Humans have long envisioned a future where humanoid robots live and work alongside us in our everyday environments. These humanoids would possess the versatility to perform diverse tasks in the environment designed for humans like Data from Star Trek: The Next Generation (1987-1994), the agility to move in complex worlds with ease like Astro from Astro Boy (2009), and the ability to collaborate with humans through natural, gentle interactions like Baymax from Big Hero 6 (2014). Achieving this vision requires endowing humanoid robots with the intelligence to move and interact as we do, expressing intent through body language, operating fluidly in human environments, and coordinating seamlessly in shared tasks. In this work, we take a step toward realizing that vision.

Humanoid robots share a similar morphology with humans, presenting a strong opportunity to acquire the same skills we have. Learning from human demonstrations provides a scalable way to develop these skills while naturally capturing human-level agility and human-like behaviors. Yet coordinating their dozens of joints and actuators poses significant challenges due to the high dimensionality of the control problem. A floating base adds another layer of complexity, as stable balancing itself remains difficult to achieve, let alone natural, human-like movement. In contrast, humans not only acquire agile skills such as sports and dancing, but are also highly versatile: we can transition through numerous motor skills as tasks demand, and solve simple unseen tasks by composing skills we have had. Endowing these robots with the same agility, naturalness, and

versatility remains a grand-challenge in robotics.

The pursuit of humanoid control has historically advanced through model-based paradigms that pair simplified dynamics with hierarchical control (*1, 2*). In practice, a few low-frequency planners generate future trajectories over coarse variables such as center of mass (CoM), momentum, contact schedule, and foot placement (*1, 3–8*). A high-frequency, low-level controller then tracks these references with a higher-fidelity model in a myopic manner (*1, 9–12*). This split keeps computation tractable, but the simplifications impose limits. Kinematic and dynamic surrogates often yield unnatural motions—such as constant CoM height or persistently bent knees (*13–15*)—that restrict robots to a small subset of their full motion range. Attempts to improve motion quality by shaping the momentum cost introduce some arm swing and upper-body rotation but remain limited to local stylistic adjustments without restoring overall human-like coordination (*16, 17*). In addition, transferring control stacks to the real world, further introduces significant unmodeled dynamics and modeling errors in highly dynamic motions (*18*), and model-based, contact-explicit planning remains difficult for contact-rich skills such as ground rolling (*19*).

To overcome these challenges, learning-based approaches, particularly reinforcement learning (RL), have emerged as promising alternatives. With appropriate reward shaping and sim-to-real transfer, RL-trained humanoids can now perform diverse locomotion behaviors, including walking on flat ground (*20*) and slopes (*21, 22*), climbing stairs (*23, 24*), running (*25*), traversing challenging terrains such as stepping stones or narrow beams (*26, 27*), and even basic compliance control (*28*). However, these successes rely on hand-crafted rewards tuned for specific tasks, which must be redesigned for each new behavior. This makes scaling up to a large repertoire of skills prohibitively costly. Moreover, hand-crafted objectives often produce unnatural movements marked by continuous stepping, bent knees, and heavy impacts, since "naturalness" and "human-likeness" cannot be easily expressed as explicit optimization terms. Although targeted fixes such as improving arm swing via centroidal momentum optimization (*29*) have been proposed, they remain case-specific and fail to generalize across diverse motions.

A promising alternative is to leverage large-scale human motion data to directly learn human-level skills. In model-based settings, libraries of human motion trajectories can serve as effective warm starts for online controllers, yielding encouraging results (*30, 31*). Yet these methods still struggle with robustness outside laboratory conditions. In comparison, RL-based learning from

3

human motion via motion tracking has produced agile and natural skills, but offers limited adaptability at deployment. Specifically, DeepMimic-style reward formulations (*32*) have demonstrated jumping and turning (*33*), single-leg balancing (*34*), and even martial-arts sequences (*35*), but these policies rely on motion-specific tuning and do not generalize beyond the motions seen during training. Improving upon it, Adversarial Motion Priors (AMPs) (*36*) learn motion "styles" for task-specific control (*37*) and generalize beyond specific motion clips, but these policies are generally not reusable across tasks and require retraining from scratch.

To improve versatility, prior work has taken two main routes. The first is hierarchical control, which combines a task-agnostic motion tracker (*38–41*) with a task-level motion planner (*42, 43*). This enables reuse of learned skills for new tasks but typically sacrifices agility and naturalness, and, due to decoupled training, suffers from planner–controller mismatches (*44, 45*). To avoid this mismatch, the second route uses multi-task generative models to learn the motion distributions directly. Among these, variational autoencoder (VAE)–based models have demonstrated agility and skill composition (*46, 47*). However, their versatility remains limited as they depend on explicit goal conditioning during training and generalize poorly to tasks with implicit or hard-to-specify objectives, such as obstacle avoidance or long-horizon navigation, yielding out-of-distribution, jittery motions and degraded naturalness (*48*).

Here, we present BeyondMimic, a scalable framework that enables humanoid robots to achieve human-level agility, naturalness, and versatility by learning from diverse, unlabeled human motions and synthesizing them seamlessly to solve unseen real-world tasks via online optimization at inference time. Shown in Fig. 1 and Movie S1, our approach deploys outdoors across various conditions, demonstrating agile and natural motions that surpass prior work in humanoid robotics. Furthermore, we present the first real-world deployment of a flexible unified controller that handles diverse downstream tasks with a variety of skills through online planning alone, without any task-specific training, fine-tuning, or policy optimization.

BeyondMimic is built on two key insights. First, we depart from the common practice of introducing complex RL formulation to foster successful sim-to-real transfer. In conventional setups, heavy domain randomization, reward regularization, and complex observations are used to compensate for idealized simulation dynamics during real-world deployment. However, this often degrades performance, necessitating motion-specific tuning. Instead, we demonstrate that

4

scalable, high-quality motion tracking can emerge from a compact, principled formulation. To achieve a generalizable formulation without any motion-specific tuning, we carefully model the robot actuation based on classical mechanics principles and ensure proper system implementation to minimize deployment discrepancies, such as delays. This disciplined design allows us to apply domain randomization only to truly uncertain physical properties, maintaining robustness without diluting the control objective. As a result, we are able to distill our reward formulation down to only three regularization terms essential for physically consistent behaviors, in addition to a unified task reward that generalizes across motions. With a simple yet effient formulation, RL training becomes straightforward: we use the same model parameters, reward function, and hyperparameters across a diverse set of skills and different physics engines, allowing human-like behaviors to transfer zero-shot to real humanoids.

Second, versatility must extend beyond training-time diversity. The robot must adapt at test time to novel objectives after deployment. BeyondMimic synthesizes the atomic skills trained by RL into novel sequences for zero-shot, task-specific control at test time. The key concept is that diffusion models, unlike VAEs or AMPs, not only capture the complex, multimodal distribution of diverse skills, but also support online optimization for new objectives at test time. Specifically, they learn the gradient fields of the data distribution (*49*) rather than the distribution itself, enabling gradient-based optimization toward arbitrary differential objectives at test time, a technique called *classifier guidance*. To leverage this capability, we employ a state-action co-diffusion model that acts in a predictive control manner, allowing cost formulation on both future states and actions. This distinct design choice enables a single, unified controller to solve a wide range of unseen tasks with unlabeled data, such as waypoint navigation, joystick teleoperation, and obstacle avoidance, while preserving the natural style and dynamic quality of the original human motions.

In summary, BeyondMimic offers the first unified framework to achieve human-level agility, human-like behaviors, and zero-shot task versatility. By combining scalable human motion learning with diffusion-based online optimization, it bridges the gap between specialized motion tracking and general-purpose task adaptability. This work provides a practical foundation for humanoid robots capable of learning, adapting, and composing skills to move and act seamlessly within human environments.

# Results

Our framework successfully learns to control a humanoid robot to perform a wide repertoire of complex, human-like behaviors. We organize our results into two parts. First, we demonstrate scalable learning of diverse human motions, highlighting both the diversity of acquired skills and their human-like quality, including agility and naturalness. Second, we show that these learned skills can be flexibly composed and repurposed at inference time. Using diffusion-based online optimization, our controller adapts to new, unseen tasks such as command-conditioned locomotion, scene-aware navigation, and motion completion from keyframes, while enabling smooth transitions between skills without retraining.

## Scalable Learning from Human Motions

Our framework demonstrated scalable learning from diverse human motions. With motion tracking, the robot successfully learned a broad set of agile, stylish, human-like skills using a *single* formulation and *shared* hyperparameters, without motion-specific tuning.

### Diverse Skills

In total, we trained on approximately 2.5 hours of diverse human motions and successfully validated all motions in high-fidelity simulation. To assess sim-to-real transfer, we deployed 30 representative clips (totaling 15 minutes) on the physical robot, showing that the skills learned in simulation transferred reliably to hardware (Fig. 2). The demonstrated behaviors spanned static and balance-critical motions such as single-leg standing and standing up varying poses; highly dynamic skills such as single-leg jumps, turn kicks, forward jumps with 180° and 360° spins, and cartwheels; and stylized or expressive behaviors such as elderly-style walking, dance sequences, and sport moves, with recognizable human-like quality. Importantly, many challenging clips were trained jointly with multiple other skills, with the entire reference motion exceeding three minutes. Despite this, the policy is able to retain agility and stylistic details, indicating that our unified formulation scales across diverse human motions without motion-specific tuning. The list of validated motions is provided in Movie S2.

## Human-level Agility

We evaluated agility in outdoor environments with soft soil, decaying leaves, and uneven ground, introducing deformable and unstable contacts that were not present during training. Despite these unseen challenges, the robot completed complex acrobatic sequences and martial-arts-inspired motions (Fig. 3), including aerial cartwheels, a skill that requires explosive takeoff, controlled mid-air rotation, and accurate landing, and is difficult even for trained adult humans. Specifically, during the airborne phase, the robot reached a peak acceleration of $31 \, \text{m/s}^2$ and a pelvic angular velocity of up to $20 \, \text{rad/s}$ (mean $7.01 \, \text{rad/s}$). We noted that comparable values have been reported in skilled human aerial motions ($7.75 \, \text{rad/s}$ on average) (*50*). Furthermore, the landings were clean and required minimal recovery, indicating accurate airborne posture control. These results demonstrated that our framework achieved human-level agility in highly dynamic motions.

Our framework further exhibited human-like agility in contact-rich control. This includes two consecutive cartwheels, crawling on the ground, and jumping up from the ground. It further reproduced expressive athletic gestures that require coordinated whole-body timing, such as Cristiano Ronaldo's celebration jump-turns, and was able to repeat the celebration five times in succession without loss of stability or style, whereas prior work (*33*) reported only a single execution.

## Natural, Human-like Behaviors

Our framework demonstrated human-like natural behaviors across diverse tasks such as walking and running, as shown in Fig. 4A-B. To assess biomechanical similarity, we compared Ground Reaction Force (GRF) profiles from the robot and humans using force-sensing treadmills (*51, 52*), normalized by total body weight. As shown in Fig. 4C, the robot exhibited comparable GRF shapes, including double peaks with a distinct heel-strike peak and a propulsive push-off peak in walking, and single peaks in running, with aligned loading and push-off timing. In walking, the robot showed sharper force peaks, reflecting reduced compliance and leg-spring behavior. We attributed this to the lack of a toe joint, which limited rollover and push-off during stance. In running, this limitation was less pronounced because contact occurred only once per step and for a shorter duration, reducing the need for a compliant toe-off phase.

Furthermore, to assess perceptual naturalness quantitatively, we conducted a user study ($N = 77$)

comparing our motions with Unitree's native controller, which is the state of the art for this robot. Participants viewed 20 paired 5-second clips of walking and running from our framework and from Unitree's native controller, selecting which motion appeared "more human-like and natural." We conducted a two-tailed binomial test to examine overall preference ($\alpha = .05$), with separate tests for each gait type using Bonferroni correction ($\alpha = .025$). BeyondMimic was strongly preferred overall (70.8% vs. 29.2%, $p < .001$, Cohen's $h = 0.859$), with significant preferences in both walking (57.0% vs. 43.0%, $p < .001$, $h = 0.281$) and running (84.7% vs. 15.3%, $p < .001$, $h = 1.532$).

Finally, we demonstrated human-like recovery under external disturbances (Fig. 4F). When the robot was gently held while walking, it remained compliant, paused and stabilized in place, and then smoothly resumed walking once released, avoiding stiff or exaggerated reactions.

## Versatile Humanoid Control

A key advantage of our framework is its ability to synthesize agile, human-like skills at inference time via diffusion guidance, enabling the completion of unseen tasks without retraining. Figure 5A illustrates this predictive control process, where the model first predicts future states and actions, then iteratively refines its predictions and converges toward a trajectory that minimizes the specified velocity cost. Importantly, this approach differs from online trajectory optimization: because the model has already acquired a diverse and feasible set of motor skills as a prior, simple, task-specific costs are sufficient to trigger appropriate behaviors for unseen tasks. This eliminates the need for numerous regularizations and behavior-shaping terms required in either model-based or learning-based methods, providing a much more versatile solution. Our experiments demonstrate this versatility in three aspects: command-conditioned locomotion with velocity and waypoint commands, inpainting using keyframes for agile skill composition, and flexible task composition (Fig. 5, Fig. 6, and Movie S5, Movie S6).

### Command-conditioned Locomotion

First, we demonstrated zero-shot commanded locomotion using either desired velocities or waypoint goals. In this task, the policy received joystick commands specifying linear and yaw velocities, or waypoint targets indicating desired positions (Fig. 5). Under waypoint navigation (Fig. 5B), it

generated smooth and stable trajectories to reach the goal from varying initial positions. Under joystick control (Fig. 5C), the robot exhibited smooth omnidirectional walking and reliably tracked the commanded direction. The controller remained robust to large external perturbations, such as kicks, while continuing to pursue the task objectives. For a longer-horizon evaluation, the robot was able to run continuously for over 50 m on a track, indicating stable control over extended distances. We recorded an average velocity tracking error of 12.14% and 13.65% in walking and running when evaluated in simulation.

The policy also demonstrated multimodal locomotion behaviors, producing distinct gaits under similar commands. For example, a low-velocity command led to either a stable walking gait or a light jogging gait, both human-like and dynamically stable. Moreover, given only a desired velocity input, the policy achieved smooth gait transitions from walking to running (Fig. 5D-E). Notably, such transitions were rare and unlabeled in the motion data; yet our controller inferred and reproduced them naturally when the task required it, much like humans acquiring smooth skill transitions through context and intention rather than explicit skill specification.

**Motion Inpainting and Task Transitioning**

Having established multimodal locomotion under velocity commands and waypoints, many agile acrobatics demand stronger, temporally structured objectives. We therefore move beyond weak velocity or position cues and instead use motion inpainting, where a sparse set of future keyframes guides the policy to generate smooth intermediate motions, enabling online insertion, transition, and composition of agile motor skills. As illustrated in Fig. 6A, starting from joystick-controlled walking, we injected desired cartwheel keyframes at 0.2 s intervals. Consequently, the diffusion policy drove a smooth transition into the cartwheel and inpainted the discrete keyframes into a temporally coherent and continuous trajectory. Upon completion of the cartwheel, the controller returned to command conditioned walking smoothly.

Following the same inpainting setup, the diffusion model was able to accurately demonstrate a diverse range of challenging motions learned from human motion tracking, along with the emergent transitions into these motions. Starting from command-conditioned locomotion, the model transitioned smoothly into these complex skills, including contact-rich behaviors such as walking into lying down and agilely getting back up to standing, as well as highly dynamic gaits such as

9

spin-kicks and flip-kicks.

Beyond the versatility to compose agile motions, our unique online optimization–based framework further enabled the versatility to transition across different task specifications. we demonstrated three consecutive cartwheels stitched with walking and running before and after, as shown in Fig. 6A(iii). This long-horizon execution showcased not only the robustness and smoothness in mode switching, but more importantly the capability to switch back and forth freely between different tasks such as velocity following and motion inpainting.

**Task Composition**

Beyond transitioning, our framework demonstrated intuitive task composition, a challenge for most prior goal-conditioned controllers. This challenge arises mainly from the combinatorial growth of possible goal combinations to enumerate. Moreover, certain tasks are difficult to be fully covered during training but can be efficiently evaluated at test time. Our approach addressed this by allowing multiple objectives to be evaluated and optimized at inference time without enumerating all possibilities. Because the costs remained simple and task-specific, they could be flexibly summed up without creating a complex objective that requires extensive tuning, enabling the framework to solve compositions of unseen tasks without retraining.

As an example of versatility in task composition, we achieved simple scene-aware navigation by combining a waypoint-tracking cost with an obstacle-avoidance cost, as shown in Fig. 6B. The obstacle-avoidance cost, shaped by a signed distance field (SDF), provides distance gradients across the prediction horizon at each denoising step, steering the predicted trajectory away from obstacles. As a result, the robot successfully detoured around the obstacle and reached the target waypoint. This same cost composition remained effective when a joystick-tracking cost replaced the waypoint term, enabling the system to mitigate collisions even under moderate off-goal user inputs.

# Discussion

The presented results substantially advance the published state of the art in humanoid control. BeyondMimic is a scalable pipeline that learns directly from human motion to reproduce agile, natural behaviors and synthesize them into task-directed actions – all without task-specific tuning.

By learning from diverse human motion data, our model acquires hundreds of skills that demonstrate human-level agility while maintaining stylistic and naturalistic characteristics. Moreover, the model does not merely memorize complete motion sequences; it composes atomic behaviors seamlessly to produce novel trajectories.

Before our work, a prevailing assumption in humanoid control was that achieving robust real-world transfer required heavy domain randomization, manual gain tuning, and carefully engineered regularization to counteract the large sim-to-real gap. These efforts often involved random torque perturbations, simulated delays, or delicately tuned reward shaping terms penalizing contact forces, slippage, and stumbles. In contrast, we show that such extensive heuristics are unnecessary. With a principled reward formulation, moderate domain randomization, and careful system implementation, our framework learns hundreds of diverse motions under a single RL formulation and one shared set of hyperparameters. The same policy transfers directly to hardware—without motion- or robot-specific tuning—while retaining both agility and natural appearance.

More importantly, BeyondMimic moves beyond imitation. For the first time, we demonstrate that a unified model learned from unlabeled human motion can synthesize and compose its skills into coherent sequences that directly solve novel downstream tasks. Through guidance in the denoising process, the controller predicts and adjusts future trajectories to optimize for new objectives at inference time, effectively coupling planning and control within a single model. As a result, the same model that performs sprinting or cartwheeling can autonomously navigate around obstacles, follow joystick-driven velocity commands, or pause and recover under external disturbances, while preserving smooth, human-like behaviors.

The training process is fully task-agnostic and requires no labeled data, providing a scalable path toward general-purpose humanoid control. This design enables continuous extension: new motion data expands the skill set without requiring redefinition of objectives or retraining for each task. The open-source release of BeyondMimic has been widely adopted, serving as a default method in public reinforcement learning repositories such as MJLab and Unitree RL Lab (*53, 54*). More broadly, this work expands the domain of humanoid control from manually tuned, motion-specific policies to a predictive, generalizable paradigm that unifies high-quality human-like behaviors with control and planning. We view this as a step toward foundational models for humanoid behavior, models that learn directly from human demonstrations and generalize across motions, environments, and

objectives.

**Limtations and Future Works**

Our approach has several limitations that present opportunities for future work. The diffusion model inherits the quality of the underlying state estimation system, such that errors in proprioception directly propagate to generated trajectories. While our latent diffusion formulation provides some robustness to noisy observations, improving state estimation through sensor fusion or learned estimators remains a promising direction for enhanced performance.

Beyond perception noise, the prediction capability could also be improved. The current system predicts trajectories over a 0.64-second horizon, which is sufficient for reactive control and local obstacle avoidance but insufficient for long-horizon planning tasks that require reasoning about distant goals or obstacles requiring early anticipation. In addition, the inclusion of history is critical for stabilizing future predictions but can cause the model to become trapped in repetitive motion patterns even when guidance is active. To counteract this, we apply larger guidance weights, which can, however, destabilize the denoising process during mode switching or in high-variance states. As a result, under guided diffusion, the robot is stable once a gait orbit is established, but tends to stumble at the start and end of motions. Future work could explore removing or reducing this dependence on history to improve responsiveness and transient behaviors.

Finally, at the task level, our guidance-based optimization works well for coarse-grained objectives but less so with fine-grained ones, and still requires light-weight tuning of the guidance weights. Future work could explore established approaches for fine-grained control over diffusion models, such as supervised fine-tuning (*55, 56*) and adapter-style control layers (*57–59*), which have proven highly effective in vision domains and could potentially enable fine-grained, composable trajectory control without extensive retraining or manual weight tuning.

# Materials and Methods

## Overview

The objective of our model is to achieve versatile humanoid control on various unseen downstream tasks, synthesizing diverse motions with sustained agility and human-like naturalness.

In the first stage, we focus on learning a diverse set of human motions through scalable motion tracking with RL. Prior works either trained a single multi-skill policy (*39, 42*), which is scalable but produces unnatural behaviors due to insufficient RL exploration, or learned separate motion-specific policies (*34, 35*) that achieve natural motions but require motion-specific tuning. Contrary to these methods, we find that a general, simple yet efficient RL pipeline is sufficient for scalable motion tracking. With a shared set of hyperparameters, the pipeline preserves human-level agility and naturalness while providing the scalability and consistency needed to extend across diverse motions.

In the second stage, we train a unified diffusion model that integrates diverse motions to enable expressive skill composition and online task optimization in a predictive control manner. The key motivation for using diffusion models is that they naturally support predictive control through classifier guidance, an online optimization process that steers unconditional generation toward task-specific conditional generation. Unlike prior diffusion-based action-only policies (*60, 61*), the key to enabling this capability is the adoption of a latent state–action diffusion model, which implicitly captures how actions influence future states and provides foresight during inference. This predictive structure enables versatile control in unseen scenarios by synthesizing learned skills seamlessly.

Trainings are conducted entirely in simulation using the most accurate manufacturer-specified parameters, without any additional system identification. After training, each stage is deployed zero-shot on physical robots. A key enabler of this transfer is our deployment framework, developed entirely in C++ and optimized for real-time execution. This framework ensures stable, low-latency control and precise synchronization between policy inference and hardware execution—addressing a critical source of the sim-to-real gap. Despite running only on CPUs and a modest mobile GPU, this implementation enables strong sim-to-real transfer for both stages without any motion- or robot-specific tuning.

# Scalable Human Motion Tracking via RL

We define scalable motion tracking as a one-recipe-fits-all framework, where each motion is trained with its own policy but under a shared formulation and training setup, minimizing the tuning effort for any new motion without compromising motion quality. To achieve this goal, we present a principled, motion-agnostic pipeline for tracking human motions with human-like agility and naturalness. We formulate the motion-tracking problem as a Markov Decision Process (MDP) and solve it using RL, which maximizes the expected cumulative reward under this MDP. Given minutes-long reference motions, the pipeline produces sim-to-real-ready motion-tracking policies using the same MDP and hyperparameters across all motions. This unified setup enables seamless scaling to hundreds of skills without manual tuning.

Since our MDP formulation is time-invariant, we omit the timestep $t$ in the following formulation for clarity. The subscript 'ref' denotes quantities from the reference motion. Unless otherwise specified, all quantities in this section are expressed in the world frame.

## Tracking Objective

Given a human-retargeted reference motion, the goal is to reproduce it on a real humanoid robot with high fidelity, while tolerating global drift caused by perturbations and sim-to-real mismatch. We start from the motion data $(\mathbf{q}^{\text{ref}}, \boldsymbol{v}^{\text{ref}})$, represented as per-frame generalized positions $\mathbf{q}^{\text{ref}} = (\mathbf{p}^{\text{ref}}, R^{\text{ref}}, \boldsymbol{\theta}^{\text{ref}}) \in \mathbb{R}^3 \times \text{SO}(3) \times \mathbb{R}^{n_{\text{jnt}}}$ and velocities $\boldsymbol{v}^{\text{ref}} = (\boldsymbol{v}^{\text{ref}}, \boldsymbol{\omega}^{\text{ref}}, \dot{\boldsymbol{\theta}}^{\text{ref}}) \in \mathbb{R}^3 \times \mathbb{R}^3 \times \mathbb{R}^{n_{\text{jnt}}}$, where $n_{\text{jnt}}$ is the number of joints for the robot. Forward kinematics yields the pose $T_b^{\text{ref}} = (\mathbf{p}_b^{\text{ref}}, R_b^{\text{ref}})$ and the twist $\mathcal{V}_b^{\text{ref}} = (\mathbf{v}_b^{\text{ref}}, \omega_b^{\text{ref}})$ for each body link $b \in \mathcal{B}$, where $\mathcal{B}$ is the set of all robot bodies. To avoid redundancy from closely spaced links, a compact set of target bodies $\mathcal{B}_{\text{target}} \subset \mathcal{B}$, including the end-effectors $\mathcal{B}_{\text{ee}} \subseteq \mathcal{B}_{\text{target}}$, is selected for tracking.

Perturbations for robustness during training and the sim-to-real gap inevitably lead to global drift. To preserve motion style while allowing such drift, the policy should track relative body poses instead of global ones. As shown in Fig. 7A, we define an anchor body $b_{\text{anchor}}$ (typically the root or torso) and use it to express desired poses of the bodies in the *anchor-centered* frame. The anchor itself follows the reference directly, $T_{\text{anchor}}^{\text{des}} = T_{\text{anchor}}^{\text{ref}}$. For any non-anchor body $b \neq b_{\text{anchor}}$, we set $T_b^{\text{des}} = \mathcal{A}\left(T_b^{\text{ref}}, T_{\text{anchor}}\right)$, where $\mathcal{A}(\cdot)$ is a yaw-aligned, height-preserving transform (see Supp. S1),

while the desired twists remain unchanged $\mathcal{V}_b^{\text{des}} = \mathcal{V}_b^{\text{ref}}$. The resulting motion-tracking objective is then:

$$\left( T_{\text{anchor}}^{\text{des}}, \mathcal{V}_{\text{anchor}}^{\text{des}}, \{ T_b^{\text{des}}, \mathcal{V}_b^{\text{des}} \}_{b \in \mathcal{B}_{\text{target}}} \right).$$

This objective preserves motion style while allowing benign global drift, improving robustness and sim-to-real transfer.

**Rewards**

To maximize transferability across motions and minimize motion-specific bias, we design a simple, motion-agnostic reward composed of a unified task term and three regularization penalties. The task rewards adopt a unified, task-space formulation with uniform weighting to promote accurate body tracking across all target bodies.

We compute tracking errors for position, orientation, linear velocity, and angular velocity, $\mathbf{e}_s$ for $s \in \{\mathbf{p}, R, \mathbf{v}, \omega\}$ between the desired and actual poses and twists. For each $s$, we take the mean-squared error $\bar{e}_s$, over all target bodies $b \in \mathcal{B}_{\text{target}}$. Each $\bar{e}_s$ is mapped through a Gaussian-shaped exponential: $r(\bar{e}_s, \sigma_s) = \exp\left(-\bar{e}_s / \sigma_s^2\right)$, with nominal error $\sigma_s$ for each term determined empirically. The total task reward is defined as

$$r_{\text{task}} = \sum_{s \in \{\mathbf{p}, R, \mathbf{v}, \omega\}} r(\bar{e}_s, \sigma_s).$$

Optionally, global tracking terms can be included for the anchor body, following the same formulation as $r_{\text{task}}$ but using only the position and orientation errors, $\mathbf{e}_{\mathbf{p},\text{anchor}}$ and $\mathbf{e}_{R,\text{anchor}}$.

In contrast to prior work that uses many ad-hoc regularizations, we use only three lightweight, broadly applicable penalties that generalize across motions to encourage physically consistent behaviors. First, the joint limit penalty, $r_{\text{limit}}$, encourages joint positions to remain within soft limits to avoid damaging the hardware. Second, the action rate penalty, $r_{\text{smooth}}$, promotes smooth transitions between consecutive actions, preventing policies with excessive jitter. Third, to penalize self-collisions, we count the number of bodies whose self-contact forces exceed a predefined threshold and sum them as the total contact penalty, $r_{\text{contact}}$, over all bodies $b \notin \mathcal{B}_{\text{ee}}$. The total reward is then defined as

$$r = r_{\text{task}} - \lambda_l r_{\text{limit}} - \lambda_s r_{\text{smooth}} - \lambda_c r_{\text{contact}},$$

15

where $\lambda_l, \lambda_s, \lambda_c > 0$ denote the corresponding weights. The weights and hyperparameters are detailed in S1.

**Observation and Action**

We use a continuous, robot-centric observation space without temporal stacking, which simplifies training and improves sim-to-real transfer. The observation for policy input is

$$\mathbf{o} = [\boldsymbol{\psi}, \mathbf{e}_{\text{anchor}}, \mathcal{V}_{\text{imu}}, \boldsymbol{\theta} - \boldsymbol{\theta}^0, \dot{\boldsymbol{\theta}}, \mathbf{a}_{\text{last}}],$$

which consists of: **(i) Motion phase** from the reference motion, $\boldsymbol{\psi} = [\boldsymbol{\theta}^{\text{ref}}, \dot{\boldsymbol{\theta}}^{\text{ref}}]$, used solely as a progress cue; the policy is not intended to track these joint states directly. **(ii) Anchor pose error** $\mathbf{e}_{\text{anchor}} \in \mathbb{R}^9$, stacks the position error $\mathbf{e}_{\mathbf{p},\text{anchor}}$, a 6-D orientation error obtained by taking the first two columns of the rotation error matrix $R^{\text{des}}_{\text{anchor}} R^{\top}_{\text{anchor}}$ (Rot6D (*62*)); since the reference is defined globally, this term provides the minimal global cue needed for balance and drift correction. **(iii) Other proprioception** includes the IMU twist expressed in the IMU frame $\mathcal{V}_{\text{imu}} \in \mathbb{R}^6$, which aid push recovery, foot timing, and stabilization; joint states including joint positions relative to the default $\boldsymbol{\theta} - \boldsymbol{\theta}^0$ and joint velocities $\dot{\boldsymbol{\theta}}$; and the previous action $\mathbf{a}_{\text{last}}$. The joint state paired with the previous action offers a proxy for the joint torque and contact currently being realized. In addition, the previous action with the action-smoothness penalty $r_{\text{smooth}}$ helps suppress high-frequency jitter.

Action is designed as normalized joint position setpoints: $\boldsymbol{\theta}^{\text{sp}} = \boldsymbol{\theta}^0 + \boldsymbol{\alpha} \odot \mathbf{a}$, where $\mathbf{a} \in \mathbb{R}^{n_{\text{jnt}}}$ is the policy output and $\boldsymbol{\alpha}$ is a per-joint action scale ($\odot$ represents elementwise product). These setpoints are sent to the low-level motor drivers as position PD controller commands to generate torques. However, they are not position targets or plans to be tracked with high precision; rather, they act as intermediate variables shaping the desired torques and are intentionally not clipped by joint kinematic limits.

High joint impedance (PD gains) is common in prior work to simplify control (*32, 33, 39, 63*): in free space, the high feedback gains suppressed the natural dynamics and dominate the closed-loop dynamics, causing the joints to behave like stiff, high-bandwidth position servos—i.e., tracking becomes close to kinematic playback. However, such high-impedance settings are typically impractical on hardware: they amplify sensor noise, decrease passive compliance needed for impact absorption, and obscure implicit torque information carried by current joint state and prior commands. In

Supplementary S1, we detail a heuristic procedure for selecting reasonable joint impedances and action scales on high-DoF humanoid systems; In Supplementary S2, we include ablations on joint impedance tuning.

**Domain Randomization**

Domain randomization is essential for successful sim-to-real transfer. However, excessive randomization can degrade policy performance, making behaviors overly conservative and training unnecessarily difficult. To enable scalable training across diverse motions while preserving smooth and natural behavior, we find that maintaining a compact set of domain randomizations is most effective. In practice, we randomize three key parameters that generalize well across motions: the ground friction and restitution coefficients, the default joint positions $\theta^0$ (for both actions and observations, simulating joint offset calibration errors), and the torso center of mass position. Additionally, random velocity perturbations are applied during training to improve robustness to environmental variability. The details for domain randomizations are in S1.

**Adaptive Sampling**

Training long motion sequences presents the challenge that different segments vary widely in difficulty. Uniformly sampling across the entire trajectory, as commonly done in prior works (*33,63*), often oversamples easy portions while undersampling harder ones, leading to slow learning or even failure to converge. To address this, we employ an adaptive sampling strategy that prioritizes segments with higher empirical failure rates. After the policy masters the difficult segments and the failure rates decrease, the sampler gradually reverts to uniform sampling to preserve coverage of easier regions.

## Versatile Humanoid Control via Guided Diffusion

The key to achieving versatile control is enabling online optimization for tasks unseen during training. Diffusion models (*64*), beyond capturing multimodal distributions, can optimize their outputs toward novel objectives via gradient ascent, a process known as classifier guidance. First introduced for image generation (*65*), it has been proven effective in RL (*66*) and character anima-

tion (*67, 68*). However, using this capability is nontrivial, as task objectives are defined in the state space, whereas typical policies output in the action space. One approach is to use forward dynamics models (*69, 70*) to bridge this gap, but such models are often constrained by high dimensionality and real-time requirements. Inspired by prior works, we instead bridge this gap by jointly modeling states and actions in a latent diffusion model with causal consistency between them. Moreover, instead of modeling only the next step, the model predicts a short horizon of future trajectories, allowing task objectives to steer the future states, which in turn produce the corresponding actions in a receding-horizon control fashion.

Shown in Fig. 7B, the latent state–action diffusion model consists of two parts. First, we train a Variational Autoencoder (VAE) using a diverse set of motion-tracking policies, which provides smooth motion representations for the diffusion model. Second, we train a state–latent diffusion model using latents collected by rolling out the VAE. Both components are trained on unlabeled, task-agnostic data. At inference time, we apply an unseen cost function to steer the state–latent diffusion model toward desired trajectories and decode its latent outputs with the VAE decoder to produce joint-level actions.

**Latent Diffusion Models**

Diffusion models are a class of generative models that generate samples by reversing a gradual noising process (*64*). In the forward process, a clean sample $\mathbf{x}_0$ is progressively corrupted by Gaussian noise, and in the reverse process, pure Gaussian noise is iteratively denoised using Stochastic Langevin Dynamics (*71*) to produce a clean sample.

Building upon this framework, Latent Diffusion Models (LDMs) (*72*) perform diffusion in a lower-dimensional latent space, typically obtained via a VAE with encoder $\mathcal{E}$ and decoder $\mathcal{D}$ that map between data and latent spaces as $\mathbf{z} = \mathcal{E}(\mathbf{x})$ and $\hat{\mathbf{x}} = \mathcal{D}(\mathbf{z})$. The VAE is trained by maximizing the evidence lower bound (ELBO):

$$\mathcal{L}_{\text{VAE}} = \mathbb{E}_{q_{\mathcal{E}}(\mathbf{z}|\mathbf{x})} \left[ \|\mathbf{x} - \mathcal{D}(\mathbf{z})\|^2 \right] + \beta \, D_{\text{KL}} \left( q_{\mathcal{E}}(\mathbf{z}|\mathbf{x}) \,\|\, \mathcal{N}(\mathbf{0}, \mathbf{I}) \right),$$

where $D_{\text{KL}}$ denotes the Kullback–Leibler divergence (*73*), and $q_{\mathcal{E}}(\mathbf{z}|\mathbf{x})$ is the encoder's posterior distribution. The first term promotes accurate reconstruction, and the second, scaled by $\beta$, regularizes the latent space.

The LDM used in this work follows the denoising diffusion probabilistic model (DDPM) (*64*), which defines a noise schedule parameterized by coefficients $\{\alpha_k, \gamma_k, \sigma_k\}_{k=1}^{K}$, controlling the signal retention and noise magnitude up to the maximum diffusion step $K$. Let $\mathbf{z}^k$ denote the noisy latent at diffusion step $k$. In the forward diffusion process, the clean data $\mathbf{x}^0$ is first encoded into its latent representation $\mathbf{z}^0 = \mathcal{E}(\mathbf{x}^0)$, and then noised to step $k$ following the diffusion posterior distribution,

$$q_{\text{forward}}(\mathbf{z}^k \mid \mathbf{z}^0) = \mathcal{N}\left(\sqrt{\bar{\alpha}_k}\,\mathbf{z}^0,\ (1 - \bar{\alpha}_k)\mathbf{I}\right),$$

where $\bar{\alpha}_k = \prod_{i=1}^{k} \alpha_i$. Then, a neural network $z_\phi(\mathbf{z}^k, k)$ is trained to predict the clean latent from the noisy input via

$$\mathcal{L}_{\text{Diffusion}} = \mathbb{E}\left[\|z_\phi(\mathbf{z}^k, k) - \mathbf{z}^0\|^2\right].$$

During inference, the reverse process iteratively denoises Gaussian noise back into a clean latent:

$$\mathbf{z}^{k-1} = \alpha_k\left(\mathbf{z}^k - \gamma_k\left(\mathbf{z}^k - z_\phi(\mathbf{z}^k, k)\right)\right) + \sigma_k\,\mathcal{N}(0, \mathbf{I}).$$

Finally, the decoder $\mathcal{D}$ maps the recovered latent $\mathbf{z}_0$ back to the data space as $\hat{\mathbf{x}}^0 = \mathcal{D}(\mathbf{z}^0)$.

## LDMs for Trajectory Modeling

In our setting, predictive control with human-like behavior amounts to learning the distribution of trajectories that give rise to coordinated, human-like motion. We formulate this as a trajectory modeling problem and use an LDM to represent it. The trajectories are generated by motion-tracking policies and serve as samples of general human-like behavior. Importantly, they include only state-action pairs, excluding any reference information, since the goal is to model the intrinsic distribution of human-like behaviors rather than the tracking process itself.

To ensure spatial consistency, we use a hybrid character–yaw-centric parameterization for the state space. A trajectory consists of $N$ past timesteps, current timestep, and $H$ future timesteps. For each timestep $n \in [-N, H]$ in the trajectory, the state $\mathbf{s}_{t+n}$ contains root pose $(\mathbf{p}_{\text{imu}}^{t+n}, \mathbf{R}_{\text{imu}}^{t+n})$ and twist $(\mathbf{v}_{\text{imu}}^{t+n}, \boldsymbol{\omega}_{\text{imu}}^{t+n})$ expressed with respect to the current root frame at timestep $t$, and body positions $\mathbf{p}_{\text{b}}^{t+n}$ and velocities $\mathbf{v}_{\text{b}}^{t+n}$ for bodies $b \in \mathcal{B}_{\text{target}}$ expressed relative to their local root frame at timestep $t+n$. The action space is the same as in motion-tracking policies, which correspond to the target inputs of the impedance controller. Details for state parameterization are provided in Section S2.

Given the trajectory modeling problem, the motivations to use an LDM over a standard diffusion model are twofold. First, the action space, serving as setpoints for the impedance controller, is highly irregular with sharp torque spikes, making direct diffusion learning unstable and violating the smoothness assumption of diffusion models. Second, large diffusion networks required for accurate trajectory modeling introduce non-negligible inference latency, causing generated actions to lag behind the latest states. In contrast, a latent space with smooth motion representations is well suited for diffusion modeling, and a lightweight decoder can leverage up-to-date observations to transform them into precise, dynamically consistent actions.

Learning an LDM from human motions consists of two stages. First, we acquire motion representations in a latent space by imitating motion tracking policies that receive observation $\mathbf{o}$ and output action $\mathbf{a}$. Inspired by prior work (74), we encode reference motions with a conditional VAE instead of raw PD actions with a regular VAE, since motion states are more structured and easier to encode than raw actions. The encoder receives only the reference-motion components to produce a latent representation $\mathbf{z} = \mathcal{E}(\boldsymbol{\psi}, \mathbf{e}_{\text{anchor}})$ that captures motion intents. The decoder then combines this latent with other proprioceptive inputs to reconstruct the action,

$$\hat{\mathbf{a}} = \mathcal{D}(\mathbf{z}, [\mathbf{g}, \mathcal{V}_{\text{imu}}, \boldsymbol{\theta}, \dot{\boldsymbol{\theta}}, \mathbf{a}_{\text{last}}]),$$

where $\mathbf{g}$ denotes the projected gravity vector expressed in the root frame. Subsequently, the VAE is trained with DAgger (75) using the modified ELBO:

$$\mathcal{L}_{\text{VAE}} = \mathbb{E}_{q_{\mathcal{E}}(\mathbf{z}|[\mathbf{c}, \mathbf{e}_{\text{anchor}}])}\left[\|\hat{\mathbf{a}} - \mathbf{a}\|^2\right] + \beta\, D_{\text{KL}}\big(q_{\mathcal{E}}(\mathbf{z} \mid \boldsymbol{\psi}, \mathbf{e}_{\text{anchor}}) \,\|\, \mathcal{N}(\mathbf{0}, \mathbf{I})\big).$$

Second, we roll out the trained VAE on human motions to collect trajectories and encode each action into a latent, yielding state–latent trajectories $\tau = [\, \mathbf{s}_{t-N}, \mathbf{z}_{t-N}, \ldots, \mathbf{s}_t, \mathbf{z}_t, \ldots, \mathbf{s}_{t+H}, \mathbf{z}_{t+H} \,]$ for LDM training. We adapt the earlier LDM formulation to this state–latent trajectory and use *individual denoising steps* for each state and latent, $\mathbf{k} = \left[\, k_{\mathbf{s}_{t-N}}, k_{\mathbf{z}_{t-N}}, \ldots, k_{\mathbf{s}_t}, k_{\mathbf{z}_t}, \ldots, k_{\mathbf{s}_{t+H}}, k_{\mathbf{z}_{t+H}} \,\right]$, enabling varying noise levels across the horizon for inpainting observations and future poses. Training for the LDM is self-supervised. Given a clean trajectory $\tau$ and uniformly sampled denoising steps $\mathbf{k}_{s_i}, \mathbf{k}_{z_i} \sim \mathcal{U}(0, K)$, the denoising network $z_\phi(\tau^{\mathbf{k}}, \mathbf{k})$ is trained to reconstruct the clean trajectory from its noised version $\tau^{\mathbf{k}}$ by minimizing

$$\mathcal{L}_{\text{Diffusion}} = \mathbb{E}\left[\|z_\phi(\tau^{\mathbf{k}}, \mathbf{k}) - \tau\|^2\right].$$

At inference, trajectories are generated by starting from Gaussian noise and iteratively denoised using

$$\tau^{\mathbf{k}-1} = \alpha_{\mathbf{k}}\Big(\tau^{\mathbf{k}} - \gamma_{\mathbf{k}}\big(\tau^{\mathbf{k}} - z_\phi(\tau^{\mathbf{k}}, \mathbf{k})\big)\Big) + \sigma_{\mathbf{k}}\,\mathcal{N}(0, \mathbf{I}),$$

which progressively denoises the sample into a clean, temporally consistent state–latent trajectory. Finally, the current action is decoded from the current denoised latent $\mathbf{z}_t$ using the most up-to-date observations. The detailed architectures and hyperparameters for the LDM are provided in Section S3.

**Online Optimization via Guidance**

To adapt the task-agnostic LDM trained above to specific tasks at inference time, we perform online optimization via *classifier guidance*. Unlike other generative models such as GANs and VAEs that learn an explicit data distribution, diffusion models learn the *gradient field* of the data distribution itself—known as the *score function* $\nabla_\tau \log p(\tau)$. Building on this property, we can transform the unconditional score function into a *conditional* one by applying Bayes' rule:

$$\nabla_\tau \log p(\tau \mid \tau^*) = \nabla_\tau \log p(\tau) + \nabla_\tau \log p(\tau^* \mid \tau),$$

where $\tau^*$ denotes the desired or optimal trajectory. This formulation enables gradient-based optimization entirely at inference time—without retraining—provided we can express the conditional gradient term $\nabla_\tau \log p(\tau^* \mid \tau)$ in a computable form. To achieve this, we approximate the conditional likelihood using a differentiable, task-specific cost function $G(\tau)$, which quantifies how well a sampled trajectory satisfies the task objective. By associating the likelihood with this cost as $p(\tau^* \mid \tau) \propto \exp(-G(\tau))$, the conditional gradient simplifies to

$$\nabla_\tau \log p(\tau^* \mid \tau) = -\nabla_\tau G(\tau).$$

This practical formulation allows the diffusion process to incorporate arbitrary differentiable cost functions as optimization objectives, enabling versatile control over diverse and unseen tasks—such as joystick control, obstacle avoidance, and motion inpainting—in an inference-time, training-free manner. Specific cost functions for each task are detailed in S3.

## Validation of the method

We present ablation studies to justify key design decisions. For motion tracking, we justify (i) notable MDP parameters, and (ii) adaptive sampling in resets.

### MDP Design and Parameters

We systematically evaluate how different MDP design choices affect sim-to-real performance. To make the evaluation sufficiently challenging and reveal meaningful differences, we select the martial arts motion shown in Fig. 8A(ii–iii), repeat each experiment three times, and record ground-truth trajectories using a Motion Capture system. The resulting local and global tracking errors are shown in Fig. 8A.

We first examine the effect of rotation representation in the observation. Baselines using quaternions show higher tracking errors, and the one with axis–angle fails once. In contrast, smooth and continuous forms such as Rot6D achieve markedly better sim-to-real performance. While prior work (62) demonstrates the benefits of Rot6D during training, our results show that the same property is important for real-world transfer, possibly because discontinuous representations induce unstable mappings that overfit to simulation-specific correlations and do not generalize on hardware.

Next, we assess the impact of adding temporal history to the observation. While prior work finds it beneficial (76), incorporating history degrades performance in our setting. We hypothesize that this results from our minimal domain randomization: additional history may encourage the policy to memorize simulation-specific state–action patterns, leading to a larger distribution shift in real-world dynamics.

We further analyze the sensitivity of sim-to-real performance to the mechanical parameters used in simulation, with armature as an example. While some practitioners modify armature values for numerical stability, we find that using inaccurate parameters noticeably degrades tracking accuracy. In particular, ignoring armature entirely (setting armature to zero) leads to excessive joint accelerations, resulting in overswinging and self-collisions during dynamic motions. More details on armature computation are presented in S1.

Finally, we investigate the effect of latency in the deployment pipeline. By injecting artificial

delays, we find that even small delays can severely impair tracking: a 2 ms delay increases velocity error, a 5 ms delay causes one failure, and a 10 ms delay causes failures in two out of three trials. While randomizing delay during training can potentially mitigate this issue, it also makes learning more difficult. These results highlight the importance of a carefully engineered, real-time deployment framework for robust sim-to-real transfer.

**Adpative Sampling**

We ablate adaptive sampling during training to evaluate its effectiveness, using convergence speed as the primary metric. As shown in Fig. 8B, adaptive sampling proves critical for learning difficult segments within long motion sequences. Without it, three out of four motions fail in certain challenging segments even after 30k training iterations—for example, the cartwheels in Motion 1 (Fig. 8C). For easier motions, such as Motion 4, adaptive sampling also halves the required iterations (2k vs. 4k). These results demonstrate that adaptive sampling substantially accelerates training and improves robustness.

To further illustrate its mechanism, we visualize the sampling process in Fig. 8D. Training begins with a uniform sampling distribution, but as difficult segments cause higher failure rates—such as the two cartwheels in Motion 1—the distribution gradually shifts to favor those segments. This reweighting reduces the variance of the policy gradient and promotes faster convergence. As the policy improves and failure rates decrease in these segments, the distribution becomes smoother and more balanced across all segments.

**Latent Diffusion**

To evaluate the contribution of our latent space formulation, we conducted an ablation study using cartwheels, a challenging athletic motion requiring precise spatiotemporal coordination. We compared success rates in sim-to-sim transfer using MuJoCo (*77*), where success was defined as completing the full cartwheel motion without falling.

The baseline method without latent encoding achieved only 5% success rate, whereas our latent diffusion model achieved 95% success in simulation. Critically, this high performance transferred to the physical robot, as shown in (Fig. 6). This improvement likely stems from enhanced robustness to variations and noise, consistent with prior work on using latent representations (*78*).

**Figure 1**: **Overview of the proposed versatile humanoid control framework.** (**A**) Scalable and robust learning from human motions with agile, human-like behaviors via motion tracking. (**B**) Versatile control over unseen downstream tasks with diverse learned motor skills via guided diffusion.

**Figure 2**: **Diverse motion tracking policies deployed on a real humanoid robot.** We demonstrate accurate and robust tracking across a wide spectrum of motions, ranging from static to highly dynamic and from athletic feats to stylized motions, showing the scalability of our framework. In total, 30 distinct motion clips are deployed on hardware; with a subset shown here due to space. See Movie S2 for the full repertoire.
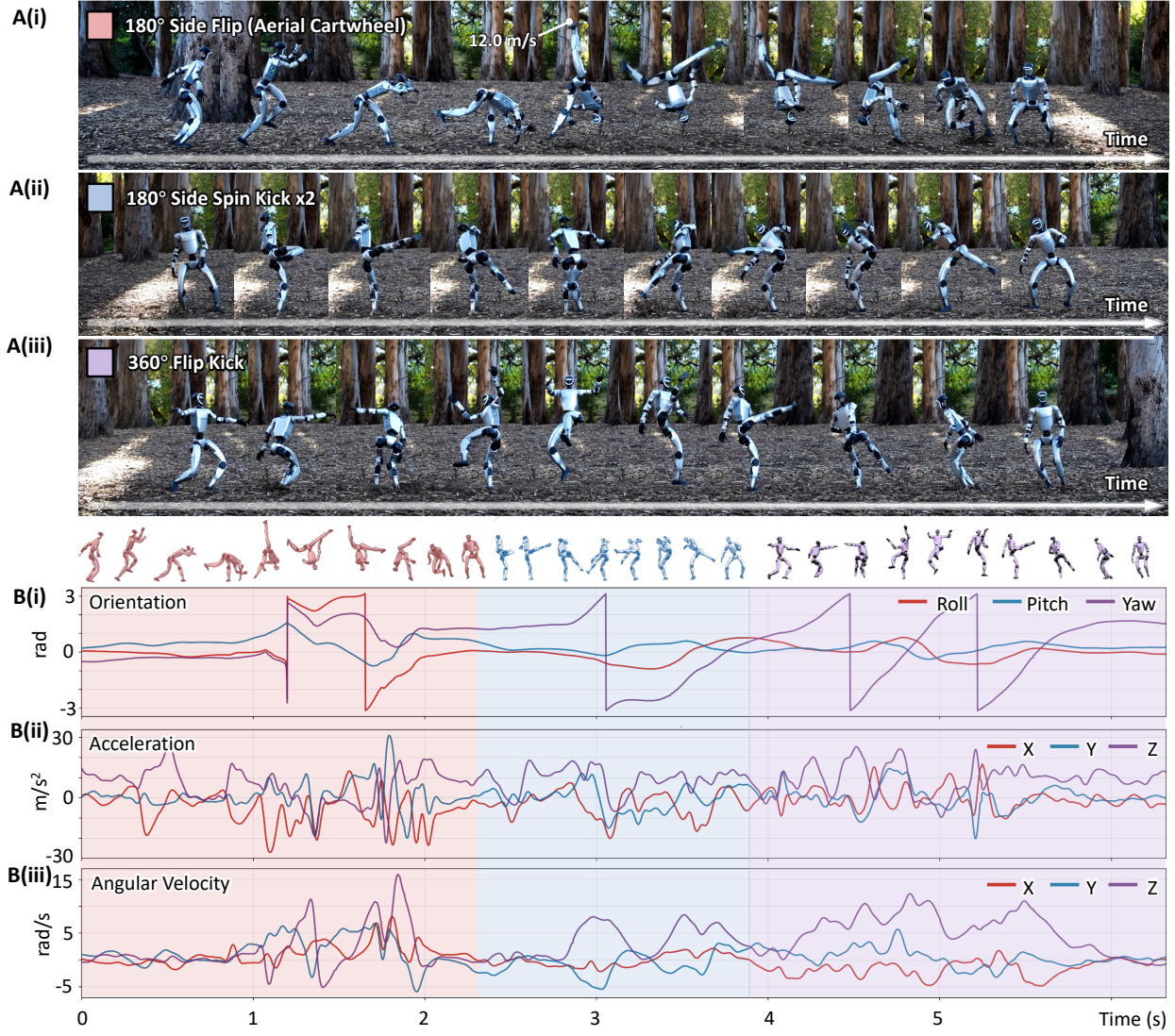
**Figure 3**: **Human-level Agility exhibited in a highly dynamic acrobatic motion.** (**A**) Real-world forest demonstration showing one continuous motion including: 180° side flip (aerial cartwheel), two consecutive 180° side spin kicks, and a 360° flip kick. (**B**) Robot's orientation, linear acceleration, and angular velocity, illustrating highly dynamic behavior; red, blue, and purple shaded intervals denote the 180° side flip, the two 180° side spin kicks, and the 360° flip kick, respectively.
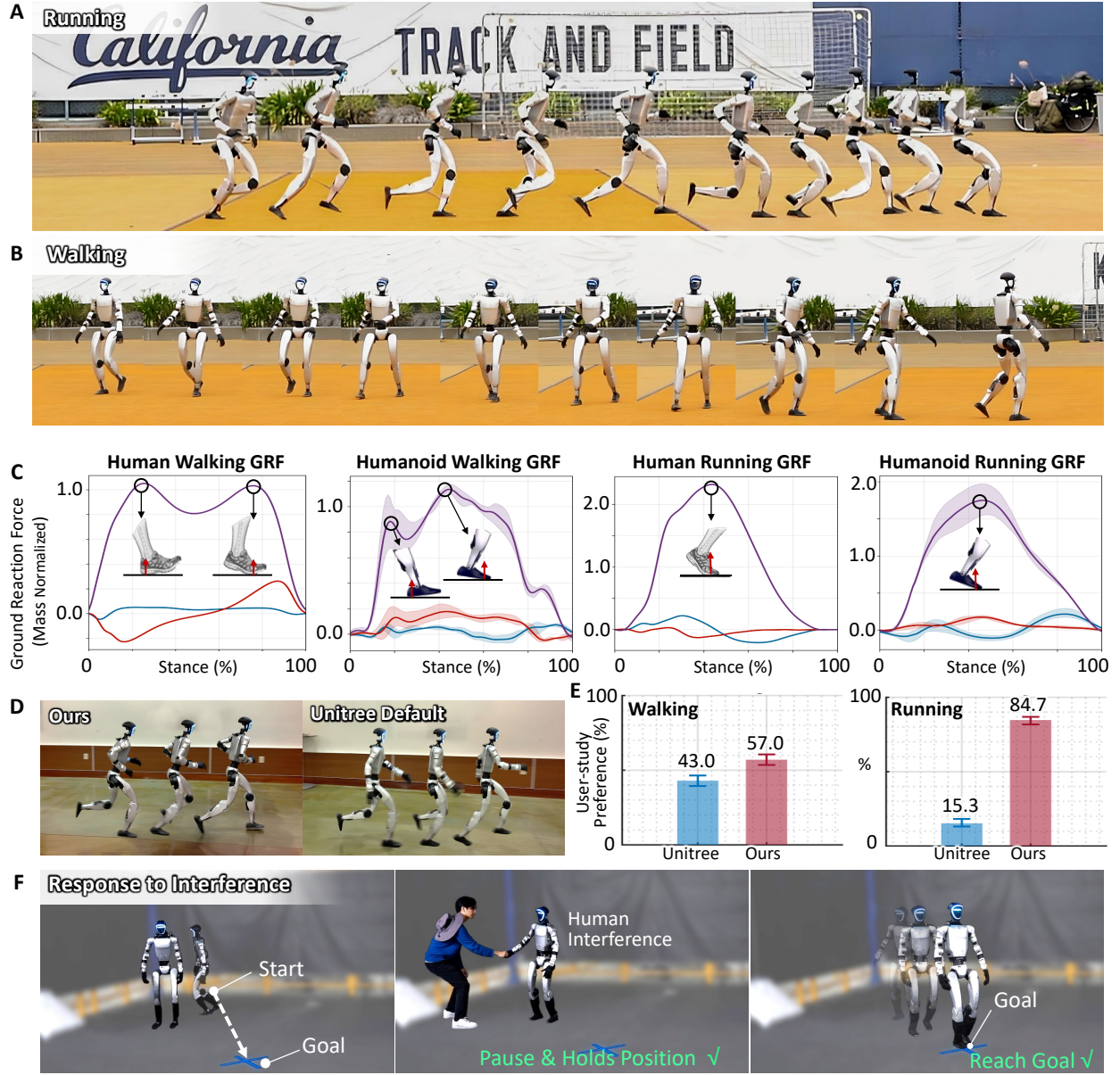
**Figure 4**: **Human-like naturalness in walking and running motion.** (**A**) Natural running on a sports field. (**B**) Natural walking on a sports field. (**C**) Comparison of weight-normalized GRF profiles for human vs. humanoid during walking and running, showing comparable shapes, contact forces, and timing. (**D**) Example video clips used in the user study. (**E**) User-study results for the question "Which looks more human-like and natural?", comparing our policies vs. Unitree's native walking and running controller, showing that ours is significantly more preferred. The error bars indicate 95% confidence intervals. (**F**) Natural response to interference during walking to a goal.
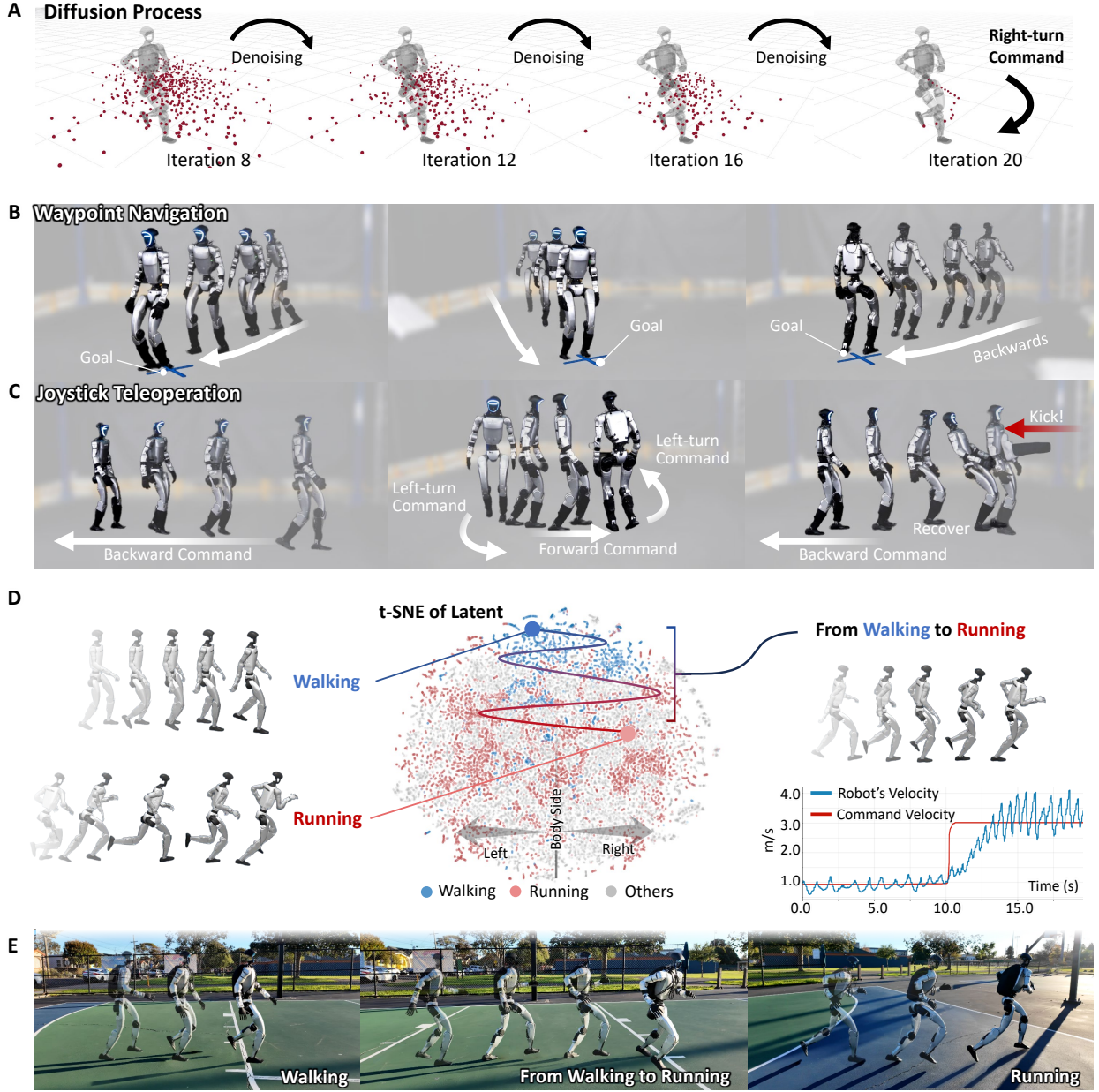
**Figure 5**: **Command-conditioned locomotion via guided diffusion.** (**A**) Visualization of diffusion process under joystick control. Starting from Gaussian noise, the distribution progressively converges to optimize for a right-turn command. (**B**) Waypoint navigation. From multiple start points, the robot reaches the goal using forward or backward walking. (**C**) Joystick teleoperation. The robot tracks the joystick velocity command. Even under an impulsive disturbance, it recovers quickly and continues following the backward command. (**D**) t-SNE visualization of the latent space, illustrating the transition from walking to running. (**E**) Real-world transition from walking to running conditioned on velocity command.

**Figure 6**: **Task transition and composition.** (**A**) Motion inpainting with future keyframes. (i) Desired future keyframes with 0.2 s intervals. (ii) Starting from walking, the robot smoothly completed the cartwheel by inpainting both the transition and the intermediate motion between the keyframes. (iii) Long-horizon execution with four keyframed-conditioned cartwheels interwoven with velocity-conditioned walking and running, achieving smooth multi-round transitions between different tasks and motions, demonstrating its versatility in task specifications. (**B**) Real-world obstacle avoidance. We demonstrated scene-aware navigation by composing waypoint and obstacle-avoidance costs. When given a goal directly ahead, the robot successfully detoured around the obstacle and reached the target.
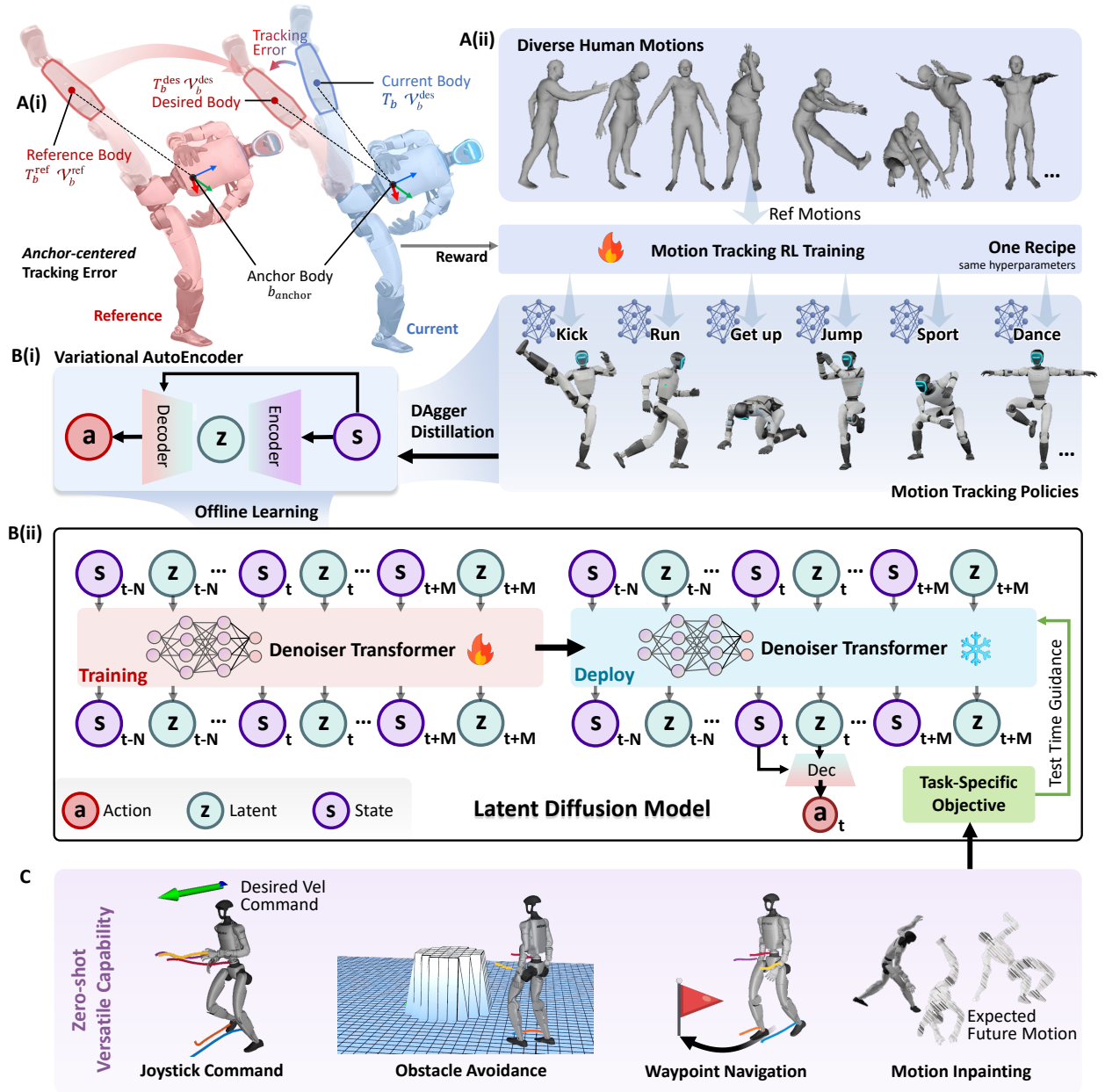
**Figure 7**: **Overview of the framework.** (**A**) Scalable learning of diverse human motions via motion tracking. (i) The target motion is re-anchored to the current pose to allow drift and recovery, which is critical for successful sim-to-real transfer. (ii) Diverse motions are tracked with RL using a single recipe and shared hyperparameters, showing scalability across a broad spectrum of skills. (**B**) Versatile control via latent state-action diffusion model. (i) Stage 1: A VAE is trained via DAgger to compress the diverse motion tracking policies into a smooth, structured latent space. (ii) Stage 2: A state-latent diffusion model is trained on VAE trajectory rollouts. During inference, the current action is decoded from the denoised latent. (**C**) Versatile capabilities demonstrated zero-shot on unseen downstream tasks through test-time guidance.
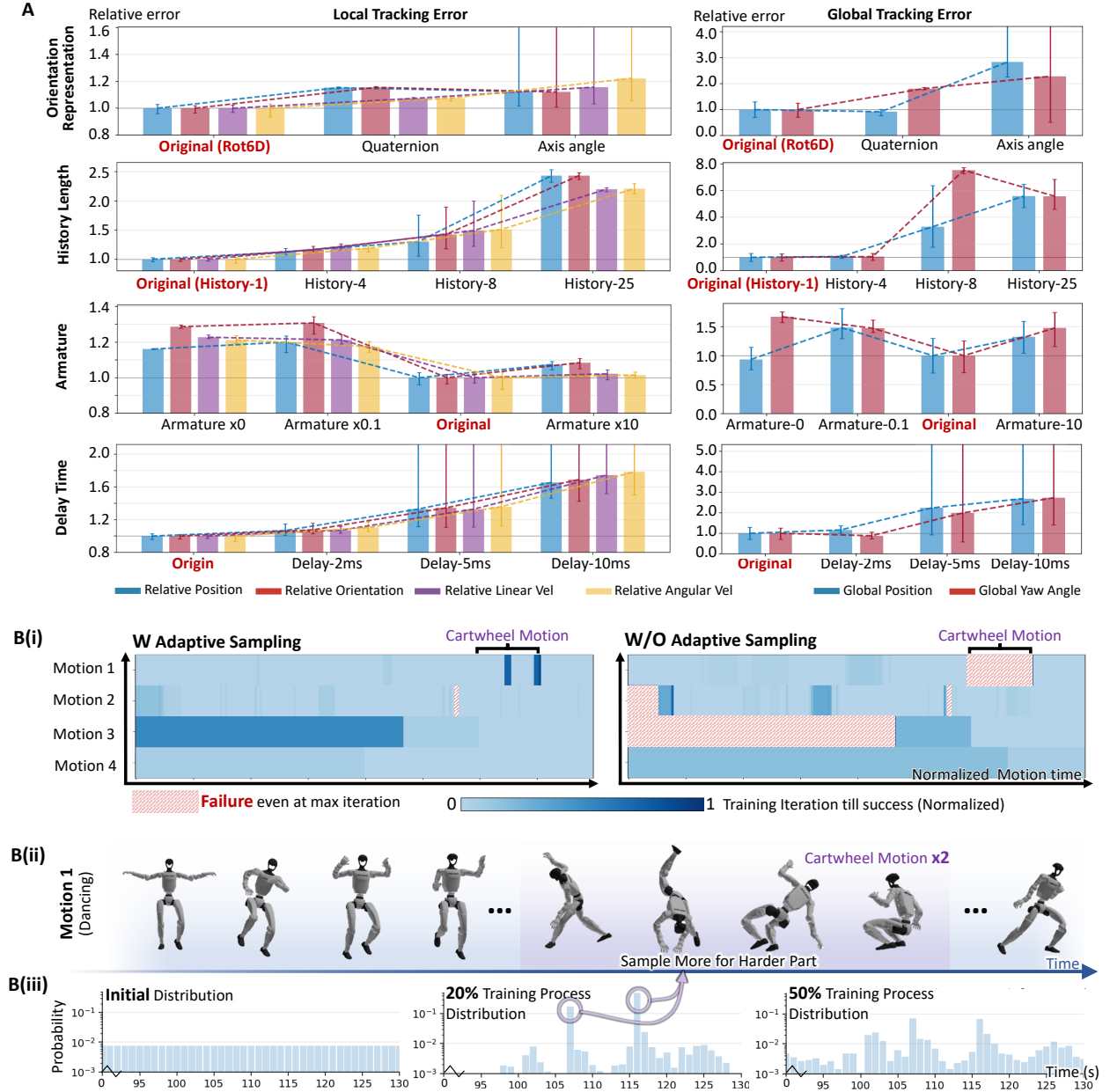
**Figure 8**: **Ablation studies.** (**A**) Motion tracking error comparison ablating orientation representation, history length, armature, and delay, showing that continuous orientations, no observation history, correct armature settings, and minimal deployment delay are critical for sim-to-real transfer in our setup. (**B**) Training performance ablating adaptive sampling (AS). (i) Iterations required to solve the motion; without AS, difficult segments remain unsolved even after 30k iterations. (ii) Visualization for *Motion 1*, highlighting two challenging cartwheel segments that cause the baseline to fail. (iii) Evolution of the sampling distribution, showing that AS concentrates sampling on difficult segments early on and spreads out as learning progresses.

# References and Notes

1. S. Kuindersma, *et al.*, Optimization-based locomotion planning, estimation, and control design for the atlas humanoid robot. *Autonomous robots* **40** (3), 429–455 (2016).

2. P. M. Wensing, *et al.*, Optimization-based control for dynamic legged robots. *IEEE Transactions on Robotics* **40**, 43–63 (2023).

3. S. Kajita, *et al.*, Biped walking pattern generation by using preview control of zero-moment point, in *2003 IEEE international conference on robotics and automation (Cat. No. 03CH37422)* (IEEE), vol. 2 (2003), pp. 1620–1626.

4. J. Pratt, J. Carff, S. Drakunov, A. Goswami, Capture point: A step toward humanoid push recovery, in *2006 6th IEEE-RAS international conference on humanoid robots* (Ieee) (2006), pp. 200–207.

5. R. Deits, R. Tedrake, Footstep planning on uneven terrain with mixed-integer convex optimization, in *2014 IEEE-RAS international conference on humanoid robots* (IEEE) (2014), pp. 279–286.

6. H. Dai, A. Valenzuela, R. Tedrake, Whole-body motion planning with centroidal dynamics and full kinematics, in *2014 IEEE-RAS International Conference on Humanoid Robots* (IEEE) (2014), pp. 295–302.

7. A. Hereid, C. M. Hubicki, E. A. Cousineau, J. W. Hurst, A. D. Ames, Hybrid zero dynamics based multiple shooting optimization with applications to robotic walking, in *2015 IEEE International Conference on Robotics and Automation (ICRA)* (IEEE) (2015), pp. 5734–5740.

8. J. Koenemann, *et al.*, Whole-body model-predictive control applied to the HRP-2 humanoid, in *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (IEEE) (2015), pp. 3346–3351.

9. 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).

10. A. Herzog, *et al.*, Momentum control with hierarchical inverse dynamics on a torque-controlled humanoid. *Autonomous Robots* **40** (3), 473–491 (2016).

11. P. M. Wensing, D. E. Orin, Generation of dynamic humanoid behaviors through task-space control with conic optimization, in *2013 IEEE International Conference on Robotics and Automation* (IEEE) (2013), pp. 3103–3109.

12. C. Khazoom, D. Gonzalez-Diaz, Y. Ding, S. Kim, Humanoid self-collision avoidance using whole-body control with control barrier functions, in *2022 IEEE-RAS 21st International Conference on Humanoid Robots (Humanoids)* (IEEE) (2022), pp. 558–565.

13. Z. Li, B. Vanderborght, N. G. Tsagarakis, D. G. Caldwell, Quasi-straightened knee walking for the humanoid robot, in *Modeling, Simulation and Optimization of Bipedal Walking* (Springer), pp. 117–130 (2013).

14. J. Carpentier, R. Budhiraja, N. Mansard, Learning Feasibility Constraints for Multicontact Locomotion of Legged Robots., in *Robotics: Science and Systems* (Cambridge, MA) (2017), p. 9p.

15. S. Fasano, J. Foster, S. Bertrand, C. DeBuys, R. Griffin, Efficient, Dynamic Locomotion through Step Placement with Straight Legs and Rolling Contacts, in *2024 IEEE International Conference on Robotics and Automation (ICRA)* (IEEE) (2024), pp. 1143–1150.

16. P. M. Wensing, 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).

17. Y.-M. Chen, G. Nelson, R. Griffin, M. Posa, J. Pratt, Integrable whole-body orientation coordinates for legged robots, in *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (IEEE) (2023), pp. 10440–10447.

18. M. Chignoli, D. Kim, E. Stanger-Jones, S. Kim, The MIT humanoid robot: Design, motion planning, and control for acrobatic behaviors, in *2020 IEEE-RAS 20th International Conference on Humanoid Robots (Humanoids)* (IEEE) (2021), pp. 1–8.

19. R. Subburaman, N. G. Tsagarakis, J. Lee, Online rolling motion generation for humanoid falls based on active energy control concepts, in *2018 IEEE-RAS 18th International Conference on Humanoid Robots (Humanoids)* (IEEE) (2018), pp. 1–7.

20. I. Radosavovic, *et al.*, Real-world humanoid locomotion with reinforcement learning. *Science Robotics* **9** (89), eadi9579 (2024).

21. I. Radosavovic, S. Kamat, T. Darrell, J. Malik, Learning humanoid locomotion over challenging terrain. *arXiv preprint arXiv:2410.03654* (2024).

22. Q. Liao, *et al.*, Berkeley humanoid: A research platform for learning-based control, in *2025 IEEE International Conference on Robotics and Automation (ICRA)* (IEEE) (2025), pp. 2897–2904.

23. J. Long, *et al.*, Learning humanoid locomotion with perceptive internal model, in *2025 IEEE International Conference on Robotics and Automation (ICRA)* (IEEE) (2025), pp. 9997–10003.

24. J. Siekmann, K. Green, J. Warila, A. Fern, J. Hurst, Blind Bipedal Stair Traversal via Sim-to-Real Reinforcement Learning. *Robotics: Science and Systems XVII* (2021).

25. Z. Olkin, K. Li, W. D. Compton, A. D. Ames, Chasing Stability: Humanoid Running via Control Lyapunov Function Guided Reinforcement Learning. *arXiv preprint arXiv:2509.19573* (2025).

26. J. He, *et al.*, Attention-based map encoding for learning generalized legged locomotion. *Science Robotics* **10** (105), eadv3604 (2025).

27. H. Wang, *et al.*, Beamdojo: Learning agile humanoid locomotion on sparse footholds. *arXiv preprint arXiv:2502.10363* (2025).

28. P. Zhi, P. Li, J. Yin, B. Jia, S. Huang, Learning a Unified Policy for Position and Force Control in Legged Loco-Manipulation, in *Conference on Robot Learning* (PMLR) (2025), pp. 652–669.

29. H. J. Lee, S. H. Jeon, S. Kim, Learning Humanoid Arm Motion via Centroidal Momentum Regularized Multi-Agent Reinforcement Learning. *IEEE Robotics and Automation Letters* (2025).

30. P. Varin, *Estimation and planning for dynamic robot behaviors* (Harvard University) (2021).

31. R. Deits, *et al.*, Robot movement and online trajectory optimization (2023), uS Patent 11,833,680.

32. X. B. Peng, P. Abbeel, S. Levine, M. van de Panne, DeepMimic: Example-guided Deep Reinforcement Learning of Physics-based Character Skills. *ACM Trans. Graph.* **37** (4), 143:1–143:14 (2018), doi:10.1145/3197517.3201311, `http://doi.acm.org/10.1145/3197517.3201311`.

33. T. He, *et al.*, Asap: Aligning simulation and real-world physics for learning agile humanoid whole-body skills. *arXiv preprint arXiv:2502.01143* (2025).

34. T. Zhang, *et al.*, HuB: Learning Extreme Humanoid Balance. *arXiv preprint arXiv:2505.07294* (2025).

35. W. Xie, *et al.*, KungfuBot: Physics-Based Humanoid Whole-Body Control for Learning Highly-Dynamic Skills. *arXiv preprint arXiv:2506.12851* (2025).

36. X. B. Peng, Z. Ma, P. Abbeel, S. Levine, A. Kanazawa, AMP: adversarial motion priors for stylized physics-based character control. *ACM Trans. Graph.* **40** (4) (2021), doi:10.1145/3450626.3459670, `https://doi.org/10.1145/3450626.3459670`.

37. H. Wang, *et al.*, PhysHSI: Towards a Real-World Generalizable and Natural Humanoid-Scene Interaction System. *arXiv preprint arXiv:2510.11072* (2025).

38. M. Ji, *et al.*, Exbody2: Advanced expressive humanoid whole-body control. *arXiv preprint arXiv:2412.13196* (2024).

39. Y. Ze, *et al.*, TWIST: Teleoperated Whole-Body Imitation System. *arXiv preprint arXiv:2505.02833* (2025).

40. Y. Li, *et al.*, CLONE: Closed-Loop Whole-Body Humanoid Teleoperation for Long-Horizon Tasks. *arXiv preprint arXiv:2506.08931* (2025).

41. K. Yin, *et al.*, UniTracker: Learning Universal Whole-Body Motion Tracker for Humanoid Robots. *arXiv preprint arXiv:2507.07356* (2025).

42. T. He, *et al.*, OmniH2O: Universal and Dexterous Human-to-Humanoid Whole-Body Teleoperation and Learning, in *Conference on Robot Learning* (PMLR) (2025), pp. 1516–1540.

43. Z. Fu, Q. Zhao, Q. Wu, G. Wetzstein, C. Finn, Humanplus: Humanoid shadowing and imitation from humans. *arXiv preprint arXiv:2406.10454* (2024).

44. M. Xu, Y. Shi, K. Yin, X. B. Peng, Parc: Physics-based augmentation with reinforcement learning for character controllers, in *Proceedings of the Special Interest Group on Computer Graphics and Interactive Techniques Conference Conference Papers* (2025), pp. 1–11.

45. X. Huang, *et al.*, Diffuse-CLoC: Guided Diffusion for Physics-based Character Look-ahead Control (2025), `https://arxiv.org/abs/2503.11801`.

46. W. Zeng, *et al.*, Behavior Foundation Model for Humanoid Robots. *arXiv preprint arXiv:2509.13780* (2025).

47. Y. Shao, *et al.*, LangWBC: Language-directed Humanoid Whole-Body Control via End-to-end Learning. *arXiv preprint arXiv:2504.21738* (2025).

48. Y. Wu, K. Karunratanakul, Z. Luo, S. Tang, UniPhys: Unified Planner and Controller with Diffusion for Flexible Physics-Based Character Control. *arXiv preprint arXiv:2504.12540* (2025).

49. Y. Song, *et al.*, Score-Based Generative Modeling through Stochastic Differential Equations, in *International Conference on Learning Representations* (2021).

50. C. Walker, J. Warmenhoven, P. J. Sinclair, S. Cobley, The application of inertial measurement units and functional principal component analysis to evaluate movement in the forward 31/2 pike somersault springboard dive. *Sports Biomechanics* **18** (2), 146–162 (2019).

51. R. Fukuchi, C. Fukuchi, M. Duarte, A public data set of running biomechanics and the effects of running speed on lower extremity kinematics and kinetics (2017), doi: 10.6084/m9.figshare.4543435.v5, `https://figshare.com/articles/dataset/A_comprehensive_public_data_set_of_running_biomechanics_and_the_effects_of_running_speed_on_lower_extremity_kinematics_and_kinetics/4543435`.

52. B. Horsak, *et al.*, GaitRec, a large-scale ground reaction force dataset of healthy and impaired gait. *Scientific data* **7** (1), 143 (2020).

53. K. Zakka, B. Yi, Q. Liao, L. Le Lay, MJLab: Isaac Lab API, powered by MuJoCo-Warp, for RL and robotics research. (2025), `https://github.com/mujocolab/mjlab`.

54. Unitree, Unitree RL Lab: Reinforcement learning implementation for Unitree robots, based on IsaacLab. (2025), `https://github.com/unitreerobotics/unitree_rl_lab`.

55. N. Ruiz, *et al.*, Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation, in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2023), pp. 22500–22510.

56. T. Brooks, A. Holynski, A. A. Efros, Instructpix2pix: Learning to follow image editing instructions, in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2023), pp. 18392–18402.

57. A. Hertz, *et al.*, Prompt-to-Prompt Image Editing with Cross-Attention Control, in *The Eleventh International Conference on Learning Representations*.

58. L. Zhang, A. Rao, M. Agrawala, Adding conditional control to text-to-image diffusion models, in *Proceedings of the IEEE/CVF international conference on computer vision* (2023), pp. 3836–3847.

59. C. Mou, *et al.*, T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models, in *Proceedings of the AAAI conference on artificial intelligence*, vol. 38 (2024), pp. 4296–4304.

60. Y. Wang, *et al.*, HDC: Humanoid Diffusion Controller, `https://humanoid-diffusion-controller.github.io/` (2025), project website with draft manuscript.

61. X. Huang, *et al.*, DiffuseLoco: Real-Time Legged Locomotion Control with Diffusion from Offline Datasets, in *8th Annual Conference on Robot Learning* (2024), `https://openreview.net/forum?id=nVJm2RdPDu`.

62. Y. Zhou, C. Barnes, J. Lu, J. Yang, H. Li, On the continuity of rotation representations in neural networks, in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2019), pp. 5745–5753.

63. Z. Luo, J. Cao, A. Winkler, K. Kitani, W. Xu, Perpetual Humanoid Control for Real-time Simulated Avatars, in *2023 IEEE/CVF International Conference on Computer Vision (ICCV)* (2023), pp. 10861–10870, doi:10.1109/ICCV51070.2023.01000.

64. J. Ho, A. Jain, P. Abbeel, Denoising diffusion probabilistic models. *Advances in neural information processing systems* **33**, 6840–6851 (2020).

65. A. Nichol, *et al.*, Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741* (2021).

66. M. Janner, Y. Du, J. B. Tenenbaum, S. Levine, Planning with diffusion for flexible behavior synthesis. *arXiv preprint arXiv:2205.09991* (2022).

67. K. Karunratanakul, K. Preechakul, S. Suwajanakorn, S. Tang, Guided motion diffusion for controllable human motion synthesis, in *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 2151–2162.

68. S. Cohan, G. Tevet, D. Reda, X. B. Peng, M. van de Panne, Flexible motion in-betweening with diffusion models, in *ACM SIGGRAPH 2024 Conference Papers* (2024), pp. 1–9.

69. H. Xue, C. Pan, Z. Yi, G. Qu, G. Shi, Full-order sampling-based mpc for torque-level locomotion control via diffusion-style annealing, in *2025 IEEE International Conference on Robotics and Automation (ICRA)* (IEEE) (2025), pp. 4974–4981.

70. P. Roth, J. Frey, C. Cadena, M. Hutter, Learned Perceptive Forward Dynamics Model for Safe and Platform-aware Robotic Navigation. *arXiv preprint arXiv:2504.19322* (2025).

71. M. Welling, Y. W. Teh, Bayesian learning via stochastic gradient Langevin dynamics, in *Proceedings of the 28th international conference on machine learning (ICML-11)* (2011), pp. 681–688.

72. R. Rombach, A. Blattmann, D. Lorenz, P. Esser, B. Ommer, High-resolution image synthesis with latent diffusion models, in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2022), pp. 10684–10695.

73. S. Kullback, R. A. Leibler, On information and sufficiency. *The annals of mathematical statistics* **22** (1), 79–86 (1951).

74. Z. Luo, *et al.*, Universal Humanoid Motion Representations for Physics-Based Control, in *The Twelfth International Conference on Learning Representations*.

75. S. Ross, G. Gordon, D. Bagnell, A reduction of imitation learning and structured prediction to no-regret online learning, in *Proceedings of the fourteenth international conference on artificial intelligence and statistics* (JMLR Workshop and Conference Proceedings) (2011), pp. 627–635.

76. J. Hwangbo, *et al.*, Learning agile and dynamic motor skills for legged robots. *Science Robotics* **4** (26), eaau5872 (2019).

77. E. Todorov, T. Erez, Y. Tassa, MuJoCo: A physics engine for model-based control, in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems* (2012), pp. 5026–5033, doi:10.1109/IROS.2012.6386109.

78. A. Serifi, R. Grandia, E. Knoop, M. Gross, M. Bächer, VMP: Versatile Motion Priors for Robustly Tracking Motion on Physical Characters. *Computer Graphics Forum* **43** (8), e15175 (2024), doi:https://doi.org/10.1111/cgf.15175.

79. Q. Liao, Dataset - BeyondMimic: From Motion Tracking to Versatile Humanoid Control via Guided Diffusion (2025), doi:10.5281/zenodo.17529720, `https://doi.org/10.5281/zenodo.17529720`.

80. V. B. Zordan, J. K. Hodgins, Motion capture-driven simulations that hit and react, in *Proceedings of the 2002 ACM SIGGRAPH/Eurographics symposium on Computer animation* (2002), pp. 89–96.

81. V. B. Zordan, A. Majkowska, B. Chiu, M. Fast, Dynamic response for motion capture animation. *ACM Transactions on Graphics (TOG)* **24** (3), 697–701 (2005).

82. M. Raibert, F. Farshidian, Workshop on Whole-body Control and Bimanual Manipulation: Applications in Humanoids and Beyond (2025), presented at the Workshop on Whole-body Control and Bimanual Manipulation: Applications in Humanoids and Beyond, Robotics: Science and Systems (RSS) 2025.

83. F. G. Harvey, M. Yurick, D. Nowrouzezahrai, C. Pal, Robust Motion In-Betweening, in *ACM Transactions on Graphics (Proceedings of ACM SIGGRAPH)* (ACM), vol. 39 (2020).

84. K. Karunratanakul, K. Preechakul, S. Suwajanakorn, S. Tang, Guided Motion Diffusion for Controllable Human Motion Synthesis, in *2023 IEEE/CVF International Conference on Computer Vision (ICCV)* (IEEE Computer Society) (2023), pp. 2151–2162, doi: 10.1109/ICCV51070.2023.00205, `https://doi.ieeecomputersociety.org/10.1109/ICCV51070.2023.00205`.

85. T. E. Truong, M. Piseno, Z. Xie, C. K. Liu, PDP: Physics-Based Character Animation via Diffusion Policy, in *SIGGRAPH Asia 2024 Conference Papers*, SA '24 (ACM) (2024), p. 1–10, doi:10.1145/3680528.3687683, `http://dx.doi.org/10.1145/3680528.3687683`.

86. G. E. Uhlenbeck, L. S. Ornstein, On the Theory of the Brownian Motion. *Phys. Rev.* **36**, 823–841 (1930), doi:10.1103/PhysRev.36.823, `https://link.aps.org/doi/10.1103/PhysRev.36.823`.

87. R. Grandia, F. Farshidian, R. Ranftl, M. Hutter, Feedback MPC for Torque-Controlled Legged Robots (2019), `https://arxiv.org/abs/1905.06144`.

88. J. Schulman, F. Wolski, P. Dhariwal, A. Radford, O. Klimov, Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* (2017).

89. T. Flayols, *et al.*, Experimental evaluation of simple estimators for humanoid robots, in *2017 IEEE-RAS 17th International Conference on Humanoid Robotics (Humanoids)* (IEEE) (2017), pp. 889–895.

90. K. Koide, M. Yokozuka, S. Oishi, A. Banno, GLIM: 3D range-inertial localization and mapping with GPU-accelerated scan matching factors. *Robotics and Autonomous Systems* **179**, 104750 (2024).

91. O. R. developers, ONNX Runtime, `https://onnxruntime.ai/` (2021), version: x.y.z.

92. A. Vaswani, *et al.*, Attention is all you need. *Advances in neural information processing systems* **30** (2017).

93. Z. Su, *et al.*, Leveraging symmetry in rl-based legged locomotion control, in *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (IEEE) (2024), pp. 6899–6906.

94. B. M. Bell, CppAD: A package for C++ Algorithmic Differentiation, `https://github.com/coin-or/CppAD/tree/stable/20190200` (2019), version 20190200.5, commit 6d82707ef4.

# Acknowledgments

**Competing interests:** There are no competing interests to declare.

**Data and materials availability:** All data needed to evaluate the conclusions in the paper are present in the paper of Supplementary Materials and in an online dataset (*79*).

## Supplementary materials

Supplementary Sections S1 to S4

Figs. S1 to S2

Tables S1 to S6

References *(33, 34, 45, 83-94)*

Movie S1 to S6

# Supplementary Materials for

# BeyondMimic: From Motion Tracking to Versatile Humanoid Control via Guided Diffusion

Qiayuan Liao[1][*][†], Takara E. Truong[2][†], Xiaoyu Huang[1][†],

Yuman Gao[1], Guy Tevet[2], Koushil Sreenath[1][‡], C. Karen Liu[2][‡]

[1]University of California, Berkeley, CA 94720, USA.

[2]Stanford University, Stanford, CA 94305, USA.

[*]Corresponding author. Email: qiayuanl@berkeley.edu

[†]These authors contributed equally to this work; order decided by coin toss.

[‡]Equal advising; order mirrors the coin toss, with the other lab listed last.

**This PDF file includes:**

Supplementary Sections S1 to S4

Figures S1 to S2

Tables S1 to S6

References *(33, 34, 45, 83-94)*

**Other Supplementary Materials for this manuscript:**

Movies S1 to S6

Dataset *(79)*.

# S1 Motion Tracking Formulation

## Armature

In a geared actuator, part of the motor's torque accelerates the motor and transmission themselves, and the remainder is transmitted through the gear train to the joint. Viewed from the joint side, the motor's rotor inertia appears as an additional inertia at the joint equal to the rotor inertia multiplied by the square of the gear ratio. Here, the gear ratio is defined as motor revolutions per joint revolution. Many simulators refer to this reflected quantity as the *armature*. Crucially, it is a real physical contribution to the joint's inertia, not a numerical stability knob; setting it incorrectly alters the robot's apparent inertia and dynamic response rather than merely the simulator's integrator behavior. Below, we detail how we compute the armature values for the four Unitree G1 robots.

As shown in Table S3, the Unitree G1 uses four two-stage planetary actuators. For completeness, we include the stage inertias in the calculation, but they contribute only a small residual–in all cases, the rotor-reflected inertia dominates the actuator armature. Building on Table S3, we map each actuator's armature to the joints summarized in Table S4. Most joints are single-actuated, so their joint armature equals the actuator armature. The exceptions are the ankle roll/pitch and waist roll/pitch joints, which are dual-actuated via a spatial linkage; in the nominal pose, this linkage is effectively 1:1, so we model their joint armature as twice the actuator value. Table S4 shows that the reflected inertia contributes a large—often dominant—fraction of the effective axis inertia in the nominal pose, especially toward the distal joints. Practically, this means controller bandwidth and gain selection at these joints are constrained primarily by the actuator's reflected inertia rather than by the link-subtree inertia.

## Anchor

For non-anchor bodies $b \in \mathcal{B} \setminus \{b_{\text{anchor}}\}$, the desired pose $T_b^{\text{des}} = (\mathbf{p}_b^{\text{des}}, R_b^{\text{des}}) = \mathcal{A}\left(T_b^{\text{ref}}, T_{\text{anchor}}\right)$ is computed as $R_b^{\text{des}} = R_\Delta R_b^{\text{ref}}, \mathbf{p}_b^{\text{des}} = \mathbf{p}_\Delta + R_\Delta(\mathbf{p}_b^{\text{ref}} - \mathbf{p}_{\text{anchor}}^{\text{ref}})$, where $\mathbf{p}_\Delta = [\mathbf{p}_{\text{anchor}.x}, \mathbf{p}_{\text{anchor}.y}, \mathbf{p}_{\text{anchor}.z}^{\text{ref}}]$ and $R_\Delta = R_z(\text{yaw}(R_{\text{anchor}} R_{\text{anchor}}^{\text{ref}}{}^\top))$. It is a transform shifts the motion into the robot's local frame by preserving height, aligning yaw, and translating the $xy$ origin under the robot.

**Reward**

To calculate the average tracking error, for each body $b \in \mathcal{B}_{\text{target}}$, we compute the position, orientation, linear and velocity tracking errors between the desired and actual poses and twists as $\mathbf{e}_{p,b} = \mathbf{p}_b^{\text{des}} - \mathbf{p}_b$, $\mathbf{e}_{R,b} = \log(R_b^{\text{des}} R_b^{\top})$, $\mathbf{e}_{v,b} = \mathbf{v}_b^{\text{des}} - \mathbf{v}_b$, and $\mathbf{e}_{\omega,b} \approx \omega_b^{\text{des}} - \omega_b$, assuming the orientation error is small. The mean squared errors are then averaged over all target bodies as $\bar{e}_s = \frac{1}{|\mathcal{B}_{\text{target}}|} \sum_{b \in \mathcal{B}_{\text{target}}} \|\mathbf{e}_{s,b}\|^2$, where $s \in \{\mathbf{p}, R, \mathbf{v}, \omega\}$. These errors are then normalized by a Gaussian-shaped exponential function: $r(\bar{e}_s, \sigma_s) = \exp\left(-\bar{e}_s/\sigma_s^2\right)$, where $\sigma$ for each term can be seen as its tolerance scale, and we choose values that fit most motions. The nominal error for each tracking term is as follows: $\sigma_p = 0.3$ for position, $\sigma_R = 0.4$ for orientation, $\sigma_v = 1.0$ for linear velocity, and $\sigma_\omega = 3.14$ for angular velocity.

The formulations for the regularization terms are as follows: $r_{\text{smooth}} = \|\mathbf{a}_t - \mathbf{a}_{t-1}\|_2$ penalizes abrupt changes in consecutive actions, $r_{\text{limit}} = \sum_{j=1}^{N_j} [\max(l_j - \theta_j, 0) + \max(\theta_j - u_j, 0)]$ penalizes joint positions $\theta_j$ exceeding the soft limits $[l_j, u_j]$, set to 0.9 times the mechanical joint limits, and $r_{\text{contact}} = \sum_{b \notin \mathcal{B}_{\text{ee}}} \mathbf{1}\left[\|\mathbf{f}_b^{\text{self}}\| > f_{\text{th}}\right]$ penalizes excessive self-contact forces, where $f_{\text{th}} = 1\text{N}$.

The weights for each term are set as follows: $\lambda_l = -10.0$, and $\lambda_s = \lambda_c = -0.1$. If the global tracking term $r_g$ is included, its weight is $\lambda_g = 0.5$. A summary of the reward terms is provided in Table S1.

**Observation**

When position drift compensation is unnecessary or reliable state estimation is unavailable, the *linear* components may be omitted (i.e., the translational part of $\mathbf{e}_{\text{anchor}}$ and the linear component of $\mathcal{V}_{\text{root}}$). We use asymmetric actor–critics to boost training efficiency. In addition to the policy observation, the critic also receives per-body relative poses w.r.t. the anchor, $T_{\text{anchor}}^{-1} T_b$, $\forall b \in \mathcal{B}$, enabling it to estimate body-wise tracking errors directly in Cartesian space.

**Heuristically Designed Paramters**

We heuristically set stiffness and damping for each joint $j$ following: $k_{\text{p},j} = I_j \omega^2$, $\quad k_{\text{d},j} = 2I_j \zeta \omega$, where $\omega$ is the natural frequency, $\zeta$ is the damping ratio, and $I_j = k_{\text{g},j}^2 I_{\text{rotor},j}$ is the reflected inertia (aramature) of the joint. Here, $I_{\text{rotor},j}$ is the rotor inertia and $k_{\text{g}}$ is the gear-ratio. The reason we

consider only the reflected inertia is that the subtree inertia varies in different configurations, and we do not rely on a heuristic to obtain a reasonable value; exact precision is not required. Note that these heuristically are built in function simulation like MuJoCo and IsaacSim for solver stability, and are used (*80–82*). We choose a damping ratio $\zeta = 2$ (overdamped) rather, since the inertia is typically underestimated by considering only the motor armature and ignoring the apparent link inertia. The natural frequency is set to a relatively low value of 10 Hz, which promotes compliance with moderate gains.

We use $\alpha = 0.25 \frac{\tau_{\max}}{k_{\mathrm{p},j}}$, with $\tau_{j,\max}$ representing the maximum allowable joint torque for joint $j$. This heuristic assumes that contacts generally occur around $\theta^0$ and that the robot hardware design ensures the maximum joint torque is proportional to the expected load.

## Domain Randomization

Here we provide the parameters for domain randomization. All ranges below denote independent uniform draws at the start of each training episode unless noted otherwise.

We randomize the robot's contact friction and restitution across all bodies, sampling static friction $\mu_{\mathrm{static}} \sim \mathcal{U}(0.3, 1.6)$, dynamic friction $\mu_{\mathrm{dynamic}} \sim \mathcal{U}(0.3, 1.2)$, and restitution $e_{\mathrm{rest}} \sim \mathcal{U}(0.0, 0.5)$. To simulate joint calibration offsets that affect both actions and observations, the default joint positions are perturbed additively by $\Delta\theta_j^0 \sim \mathcal{U}(-0.01, 0.01)$ rad for most joints and by a larger range $\Delta\theta_j^0 \sim \mathcal{U}(-0.1, 0.1)$ rad for ankle joints to account for greater calibration errors. The torso's center of mass is randomized along all three axes, with $\Delta x \sim \mathcal{U}(-0.025, 0.025)$ ref, $\Delta y \sim \mathcal{U}(-0.05, 0.05)$ ref, and $\Delta z \sim \mathcal{U}(-0.05, 0.05)$ ref. To further enhance robustness to external disturbances, random velocity perturbations are applied periodically during training, with intervals $\Delta t \sim \mathcal{U}(1.0, 3.0)$ s. The perturbation magnitudes are uniformly sampled from translational velocity ranges $(x, y, z) \in \{[-0.5, 0.5], [-0.5, 0.5], [-0.2, 0.2]\}$ m/s and rotational velocity ranges (roll, pitch, yaw) $\in \{[-0.52, 0.52], [-0.52, 0.52], [-0.78, 0.78]\}$ rad/s. A summary of the domain randomization terms is provided in Table S2.

## Termination and Adaptive Sampling

An episode terminates when the tracking error becomes unrecoverably large. The position and orientation errors are defined as $\mathbf{e}_{p,b} = \mathbf{p}_b^{\mathrm{d}} - \mathbf{p}_b$ and $\mathbf{e}_{R,b} = \log(R_b^{\mathrm{d}} R_b^{\top})$. Termination occurs when either

the anchor body $b_{\text{anchor}}$ or any end-effector body $b \in \mathcal{B}_{\text{ee}} = \{\text{left/right ankles, left/right hands}\}$ deviates excessively from the reference, that is, when $|\mathbf{e}_{p,z,b}| > 0.25$ m or $\|\mathbf{e}_{R,\text{anchor}}\| > 0.8$ rad.

At each episode reset, the motion phase is adaptively sampled from the reference trajectory to balance learning across segments of varying difficulty. The reference motion is divided into $S$ one-second bins, and each bin's failure rate is updated using an exponential moving average, $\bar{f}_s \leftarrow 0.999\,\bar{f}_s + 0.001 f_s$, where $\bar{f}_s$ is the smoothed failure rate, to mitigate short-term fluctuations. The resulting failure rates are offset by a small uniform floor, $\bar{f}_s = \bar{f}_s + 0.1/S$, so that when all bins have uniformly low failure rates, the sampler naturally reverts to uniform sampling. Since failures are more likely caused by suboptimal actions taken shortly before termination, we apply a non-causal convolution with an exponentially decaying kernel $k(u) = \rho^u$ with $\rho = 0.8$ and a look-back window of $u \in \{0, 1, 2\}$. The sampling probability of each bin s with the convolution kernel is then computed as

$$p_{\text{s}} = \frac{\sum_{u=0}^{K-1} \rho^u\,\bar{f}_{\text{s}+u}}{\sum_{j=1}^{\text{S}} \sum_{u=0}^{K-1} \rho^u\,\bar{f}_{j+u}}.$$

The probabilities are then normalized as $\hat{p}_s = p_s / \sum_{j=1}^{S} p_j$. Finally, the bin of the starting phase is drawn from $\text{Multinomial}(\hat{p}_1, \ldots, \hat{p}_S)$.

The robot is then initialized at the corresponding reference configuration and velocity for the sampled phase, with small random perturbations applied to counteract cumulative errors during long-horizon rollouts. Specifically, pose perturbations are sampled within $x, y \in [-0.05, 0.05]$ m, $z \in [-0.01, 0.01]$ m, $\text{roll}, \text{pitch} \in [-0.1, 0.1]$ rad, and yaw $\in [-0.2, 0.2]$ rad, while velocity perturbations follow the same range used in the domain randomization settings.

## S2 Motion Tracking Results

### Complete List of Validated Motions

We used various datasets to showcase the wide applicability of our framework. These include datasets from prior works (*33, 34*), Unitree-retargeted LAFAN1 (*83*), and online animation data[1]. The complete list of motions validated in simulation and on hardware is provided in Table S8.

---

[1]Motion sources: Reallusion "MD Panther Lady" pack (`https://www.reallusion.com/ContentStore/iClone/3d-animation/panther-lady/default.html`) and "Martial Arts – Taekwondo" motion pack (`https://actorcore.reallusion.com/3d-motion/pack/martial-arts-taekwondo`).

**Extra Ablation on PD Gains**

We ablate the PD gains of the impedance controller. As described earlier, the gains are parameterized by a natural frequency $\omega$, where a larger $\omega$ results in a stiffer response. We evaluate several values of $w$ and also compare against the gains used in ASAP (*33*). As shown in Fig. S2, $\omega = 10\text{Hz}$ achieves the best global tracking performance and stronger local tracking than both $\omega = 5$ Hz and ASAP. Although $\omega = 25$ Hz yields slightly lower local pose error, the higher stiffness produces noticeable high-frequency oscillation and loud sound in the actuator gearboxes, which likely results from the large torque overshoots under high gains. Therefore, to balance performance and hardware lifespan, we adopt $\omega = 10$ Hz in all experiments.

## S3 Diffusion Formulation

**Diffusion State**

We use character frames to normalize the states for consistency across trajectories. A *character frame $C_{t'}$* is defined for each timestep $t'$ by the root position $\mathbf{p}_{\text{root}}^{t'}$ and yaw rotation $\mathbf{R}_{\text{yaw}}^{t'}$ (with pitch and roll removed). All quantities expressed in $C_{t'}$ are thus invariant to global translation and yaw rotation. At the *current timestep $t$*, we define the frame $C_t$ as the current character frame for trajectory normalization. For each timestep $t + n$ within the history and horizon ($n \in [-N, H]$), the state $\mathbf{s}_{t+n}$ is consisted of root features and body features.

The root features of $\mathbf{s}_{t+n}$ are expressed relative to the current frame $C_t$, i.e., the root pose and twist are computed as

$$\mathbf{p}_{\text{root}}^{\text{rel}}(t+n) = (\mathbf{R}_{\text{yaw}}^t)^\top(\mathbf{p}_{\text{root}}^{t+n} - \mathbf{p}_{\text{root}}^t), \qquad \mathbf{R}_{\text{root}}^{\text{rel}}(t+n) = (\mathbf{R}_{\text{yaw}}^t)^\top \mathbf{R}_{\text{root}}^{t+n}, \tag{S1}$$

$$\mathbf{v}_{\text{root}}^{\text{rel}}(t+n) = (\mathbf{R}_{\text{yaw}}^t)^\top(\mathbf{v}_{\text{root}}^{t+n} - \mathbf{v}_{\text{root}}^t), \qquad \omega_{\text{root}}^{\text{rel}}(t+n) = (\mathbf{R}_{\text{yaw}}^t)^\top \omega_{\text{root}}^{t+n}. \tag{S2}$$

The body features of state $\mathbf{s}_{t+n}$ are instead expressed in their *local root frames $C_{t+n}$*. For each body $b \in \mathcal{B}_{\text{target}}$, the local positions and velocities are computed as

$$\mathbf{p}_b^{\text{local}}(t+n) = (\mathbf{R}_{\text{yaw}}^{t+n})^\top(\mathbf{p}_b^{t+n} - \mathbf{p}_{\text{root}}^{t+n}), \qquad \mathbf{v}_b^{\text{local}}(t+n) = (\mathbf{R}_{\text{yaw}}^{t+n})^\top(\mathbf{v}_b^{t+n} - \mathbf{v}_{\text{root}}^{t+n}). \tag{S3}$$

This hybrid trajectory representation expresses root motion in the current character–yaw frame and body motion in their instantaneous local frames, encoding temporal information invariant to global translation and yaw while retaining detailed local structure that eases model learning.

To enhance the temporal information in the state representation, we integrate an *emphasis projection* (*45, 84*). The projection matrix is defined as $\mathbf{P} = [\mathbf{AB}\ \mathbf{I}]^\top$, where $\mathbf{A}_{ij} \sim \mathcal{N}(0,1)$ is a random Gaussian matrix and $\mathbf{B}$ is a diagonal matrix with entries set to $c = 6$ for the dimensions corresponding to root poses and twists. This formulation again jointly encodes root features with temporal information and local body features, with the former being emphasized by the projection. The transformed input to the latent diffusion model is then $\mathbf{s}' = \mathbf{Ps}$, and after the LDM predicts the next state $\hat{\mathbf{s}}'$ in the projected space, the corresponding state in the original space is recovered via the inverse transformation $\hat{\mathbf{s}} = \mathbf{P}^{-1}\hat{\mathbf{s}}'$, where $\mathbf{P}^{-1}$ is computed using the pseudoinverse.

**Diffusion Dataset Collection**

Because the diffusion model is trained entirely on offline data, using only clean VAE rollouts often causes out-of-distribution behavior during execution. To enhance robustness, we follow PDP (*85*) to augment trajectories with policy perturbations, forming an error band that enables the model to recover from deviations during deployment. Specifically, we roll out the VAE policies with added action noise while recording the original states and latents.

Prior works (*45, 85*) add i.i.d. Gaussian action noise at each step, but the overdamped PD gains suppress such high-frequency perturbations, limiting state diversity. We instead use Ornstein-Uhlenbeck (OU) noise (*86*), which produces temporally correlated action perturbations:

$$\eta_{t+1} = \eta_t + \theta(\mu - \eta_t)\,\Delta t + \sigma\sqrt{\Delta t}\,\varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(0, I), \tag{S4}$$

where $\theta = 0.8$ controls the mean reversion rate, $\mu = 0$ is the long-term mean, and $\Delta t = 1.0$. We set a joint-wise noise scale of $\sigma = 0.1$. Actions are then perturbed as $a_t \leftarrow a_t + \eta_t$, creating a persistent *error band* that the policy learns to recover from during deployment. Episodes are collected such that each sample in the trajectory appears approximately 100 times across the dataset. Each rollout executes the policy for 2.5 seconds but allows the motion to continue for 5 seconds to verify stability. If the robot fails before 5 seconds, the episode is rejected and not included in the dataset.

**Task Costs**

We show joystick steering, waypoint navigation, and obstacle avoidance as examples of task-specific cost functions. For joystick steering, we define the cost function as the squared difference between

the state prediction and the joystick input:

$$G_{\mathrm{js}}(\hat{\tau}_t) = \frac{1}{2} \sum_{i=0}^{H} \|V_{xy,i}(\hat{\tau}_t) - \mathbf{g}_v\|^2 \tag{S5}$$

where $V_{xy,i}(\hat{\tau}_t)$ extracts the predicted planar root velocity at horizon step $i$ and $\mathbf{g}_v \in \mathbb{R}^2$ is the commanded velocity.

For waypoint navigation, the cost transitions from position tracking to velocity minimization near the goal:

$$G_{\mathrm{wp}}(\hat{\tau}_t) = \sum_{i=0}^{H} (1 - e^{-2d_i}) \|P_{xy,i}(\hat{\tau}_t) - \mathbf{g}_p\|^2 + e^{-2d_i} \|V_{xy,i}(\hat{\tau}_t)\|^2 \tag{S6}$$

where $P_{xy,i}(\hat{\tau}_t)$ extracts planar root position at horizon step $i$, $d_i = \|P_{xy,i}(\hat{\tau}_t) - \mathbf{g}_p\|$ is the distance to goal $\mathbf{g}_p \in \mathbb{R}^2$.

For obstacle avoidance, we use a Signed Distance Field (SDF) with a barrier function:

$$G_{\mathrm{sdf}}(\hat{\tau}_t) = \sum_{i=0}^{H} \sum_{b \in \mathcal{B}} B(\mathrm{SDF}(P_{b,i}(\hat{\tau}_t)) - r_b, \delta) \tag{S7}$$

where $P_{b,i}(\hat{\tau}_t)$ is the position of body $b$ at horizon step $i$, $r_b$ is its collision radius, and $B(x, \delta)$ is a relaxed barrier function (*87*):

$$B(x, \delta) = \begin{cases} -\ln(x) & \text{if } x \geq \delta \\ -\ln(\delta) + \frac{1}{2}\left[\left(\frac{x-2\delta}{\delta}\right)^2 - 1\right] & \text{if } x < \delta \end{cases} \tag{S8}$$

## S4 Implementation Details

### Motion Tracking

We use MLPs as our policy networks, and PPO (*88*) as the RL algorithm. Policies run at 50Hz. We provide the architecture and training details for motion tracking in Table S5.

We deploy our motion tracking policies, trained using the proposed method, on Unitree G1 humanoid robots. All deployment code is written in C++ and optimized for real-time execution. Full-state estimation is provided at 500 Hz using a low-level generalized momentum observer combined with a Kalman filter (*89*). Notably, no external motion capture system is used for tracking.

For extreme, contact-rich behaviors (e.g., getting up from the ground), we either incorporate LiDAR-inertial odometry (LIO) (*90*) for position correction or exclude state-estimation-dependent observations altogether. All policies are executed onboard using ONNX Runtime (*91*) on the robot's CPU. Each inference step takes under 1.0 ms, allowing seamless integration into the real-time estimation and control loop.

**Diffusion**

We use MLPs for the VAE encoder and decoder networks, and a Transformer (*92*) Encoder for the diffusion backbone. The transformer has a parameter count of $\tilde{1}9.8M$. We train and deploy tracking policies at 25 Hz instead of 50 Hz to allow sufficient time for diffusion model inference, following the exact same formulation as introduced. We provide details for the architecture and training of the VAE in Table S6 and Diffusion in Table S7.

In addition, we apply morphological symmetry augmentation in the sagittal plane when learning from the motion-tracking policies. Following the formulation in prior work (*93*), we train both the VAE and the diffusion model on two copies of the data: the original state–action pairs and their sagittal-symmetric counterparts. This effectively doubles the skill repertoire with minimal additional effort.

At test time, the diffusion policy runs onboard using a portable Mini PC with an NVIDIA RTX 4060 Mobile GPU. We use TensorRT to accelerate the inference speed. The diffusion model inference is performed asynchronously in a separate thread due to latency; each step takes approximately 20 ms using 20 denoising steps. The light-weight VAE decoder is inferenced in synchronously and on CPU instead. Gradients of the guiding cost are computed automatically using CppAD (*94*) within each denoising iteration. For the joystick commands and motion inpainting tasks, we use proprioceptive state estimation to provide body pose and velocities. For the waypoint navigation and obstacle avoidance tasks, motion capture data is used to provide both environmental context for cost computation and improved localization to accomplish the tasks.
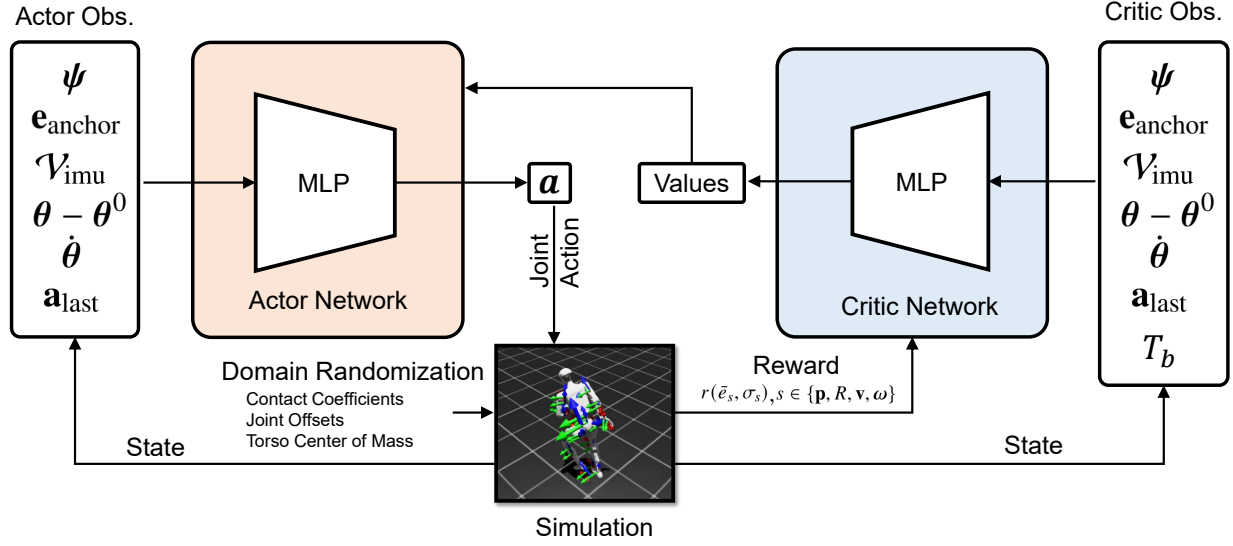
**Figure S1**: **Motion Tracking Overview.** Our actor-critic architecture for motion tracking. The actor network takes observations: motion phase $\psi$, anchor pose error $\mathbf{e}_{\text{anchor}}$, IMU twist $\mathcal{V}_{\text{imu}}$, relative joint positions $\theta - \theta^0$, joint velocities $\dot{\theta}$, and previous action $\mathbf{a}_{\text{last}}$. The critic receives the same observations in addition to the per-body relative poses w.r.t. the anchor $T_{\text{b}}$ for a desired set of bodies $b \in \mathcal{B}_{\text{target}}$. We compute tracking rewards for position $\mathbf{p}$, orientation $R$, linear velocity $\mathbf{v}$, and angular velocity $\omega$ between the desired and actual poses and twists, in addition to three regularization terms. We domain randomize over contact coefficients, joint offsets, and torso center of mass only.



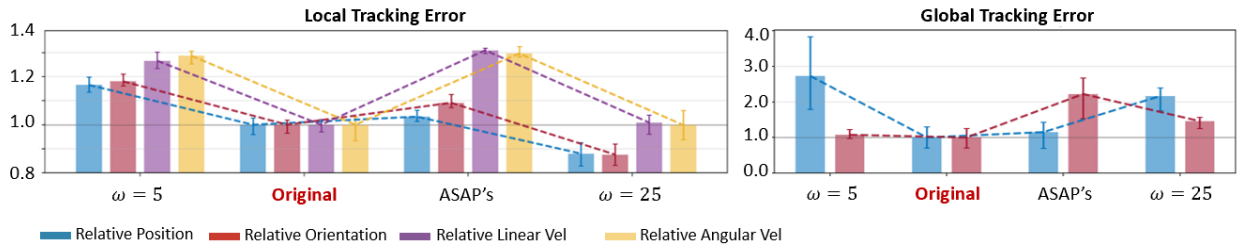**Figure S2**: **Extra Ablation On PD Gains** Ablating on PD Gains, We find that following our heuristics formula with a moderate natural frequency yields the best global tracking performance and surpasses the hand-tuned gains used in prior work (*33*). Although increasing the natural frequency can further reduce local tracking error, it results in significant torque overshoots that produces audible jitter and can stress the hardware.

**Table S1**: Unified reward formulation using Gaussian-shaped tracking scores.

| Reward Terms | Equation | Weight |
|---|---|---|
| *Task (Tracking)* | | |
| Body Position | $\exp\left(-\left(\frac{1}{|\mathcal{B}_{\text{target}}|}\sum_{b\in\mathcal{B}_{\text{target}}}\|\mathbf{p}_b^{\text{des}}-\mathbf{p}_b\|^2\right)/0.3^2\right)$ | 1.0 |
| Body Orientation | $\exp\left(-\left(\frac{1}{|\mathcal{B}_{\text{target}}|}\sum_{b\in\mathcal{B}_{\text{target}}}\|\log(R_b^{\text{des}}R_b^\top)\|^2\right)/0.4^2\right)$ | 1.0 |
| Body Linear velocity | $\exp\left(-\left(\frac{1}{|\mathcal{B}_{\text{target}}|}\sum_{b\in\mathcal{B}_{\text{target}}}\|\mathbf{v}_b^{\text{des}}-\mathbf{v}_b\|^2\right)/1.0^2\right)$ | 1.0 |
| Body Angular velocity | $\exp\left(-\left(\frac{1}{|\mathcal{B}_{\text{target}}|}\sum_{b\in\mathcal{B}_{\text{target}}}\|\boldsymbol{\omega}_b^{\text{des}}-\boldsymbol{\omega}_b\|^2\right)/3.14^2\right)$ | 1.0 |
| Anchor Position (Optional) | $\exp\left(-\|\mathbf{p}_{\text{anchor}}^{\text{des}}-\mathbf{p}_{\text{anchor}}\|^2/0.3^2\right)$ | 0.5 |
| Anchor Orientation (Optional) | $\exp\left(-\|\log(R_{\text{anchor}}^{\text{des}}R_{\text{anchor}}^\top)\|^2/0.4^2\right)$ | 0.5 |
| *Regularization* | | |
| Action smoothness | $\|\mathbf{a}_t-\mathbf{a}_{t-1}\|^2$ | −0.1 |
| Joint position limit | $\sum_{j=1}^N\left[\max(l_j-\theta_j,0)+\max(\theta_j-u_j,0)\right]$ | −10.0 |
| Undesired self-contacts | $\sum_{b\notin\mathcal{B}_{\text{ee}}}\mathbf{1}\left[\|f_b^{\text{self}}\|>1\text{N}\right]$ | −0.1 |

**Table S2**: **Domain randomization parameters.** ($\mathcal{U}[\cdot]$: uniform distribution)

| Domain Randomization | Sampling Distribution |
|---|---|
| *Physical parameters* | |
| Static friction coefficients | $\mu_{\text{static}} \sim \mathcal{U}[0.3,\ 1.6]$ |
| Dynamic friction coefficients | $\mu_{\text{dynamic}} \sim \mathcal{U}[0.3,\ 1.2]$ |
| Restitution coefficient | $e_{\text{rest}} \sim \mathcal{U}[0,\ 0.5]$ |
| Default joint positions [rad] | $\Delta\theta_j^0 \sim \mathcal{U}[-0.01,\ 0.01]$ |
| Torso's COM offset [m] | $\Delta x \sim \mathcal{U}[-0.025, 0.025],\ \ \Delta y \sim \mathcal{U}[-0.05, 0.05],$ |
| | $\Delta z \sim \mathcal{U}[-0.05, 0.05]$ |
| *Root velocity perturbations* | |
| Root linear vel [m/s] | $v_x \sim \mathcal{U}[-0.5, 0.5],\ \ v_y \sim \mathcal{U}[-0.5, 0.5],\ \ v_z \sim \mathcal{U}[-0.2, 0.2]$ |
| Push duration [s] | $\Delta t \sim \mathcal{U}[1.0,\ 3.0]$ |
| Root angular vel [rad/s] | $\omega_x,\ \omega_y \sim \mathcal{U}[-0.52, 0.52],\ \ \omega_z \sim \mathcal{U}[-0.78, 0.78]$ |

**Table S3**: **Actuator reflected inertia.** Total $J_{\text{arm}} = J_r(g_1 g_2)^2 + J_1 g_2^2 + J_2$ with unit kg·m$^2$. Near-zero entries shown as "–".

| Actuator | $g_1$ | $g_2$ | $J_r$ | $J_1$ | $J_2$ | $J_r(g_1 g_2)^2$ | $J_1 g_2^2$ | Total |
|---|---|---|---|---|---|---|---|---|
| 5020-16 | 3.56 | 4.5 | 1.390e-05 | 1.700e-06 | 1.690e-05 | 3.558e-03 | 3.443e-05 | 3.610e-03 |
| 7520-14.3 | 4.5 | 3.18 | 4.890e-05 | 9.800e-06 | 5.330e-05 | 1.003e-02 | 9.921e-05 | 1.018e-02 |
| 7520-22.5 | 4.5 | 5 | 4.890e-05 | 1.090e-05 | 7.380e-05 | 2.476e-02 | 2.725e-04 | 2.510e-02 |
| 4010-25 | 5 | 5 | 6.800e-06 | – | – | 4.250e-03 | – | 4.250e-03 |

**Table S4**: **Joint-axis inertia for left arm and left leg.** Columns list $I_{\text{axis}}^{\text{subtree}}$ (subtree inertia about the joint axis), $I_{\text{armature}}$ (reflected motor/gear inertia), $I_{\text{axis}}^{\text{eff}} = I_{\text{axis}}^{\text{subtree}} + I_{\text{armature}}$, and the percentage of armature in the effective inertia.

| Joint | $I_{\text{axis}}^{\text{subtree}}$ | $I_{\text{armature}}$ | $I_{\text{axis}}^{\text{eff}}$ | Armature % |
|---|---|---|---|---|
| | (kg·m$^2$) | (kg·m$^2$) | (kg·m$^2$) | (%) |
| *Left Arm* | | | | |
| left_shoulder_pitch_joint | 1.748e-01 | 3.610e-03 | 1.785e-01 | 2.0% |
| left_shoulder_roll_joint | 1.324e-01 | 3.610e-03 | 1.360e-01 | 2.7% |
| left_shoulder_yaw_joint | 2.742e-02 | 3.610e-03 | 3.103e-02 | 11.6% |
| left_elbow_joint | 3.445e-02 | 3.610e-03 | 3.806e-02 | 9.5% |
| left_wrist_roll_joint | 3.664e-04 | 3.610e-03 | 3.976e-03 | 90.8% |
| left_wrist_pitch_joint | 4.804e-03 | 4.250e-03 | 9.054e-03 | 46.9% |
| left_wrist_yaw_joint | 1.837e-03 | 4.250e-03 | 6.087e-03 | 69.8% |
| *Left Leg* | | | | |
| left_hip_pitch_joint | 8.543e-01 | 1.018e-02 | 8.644e-01 | 1.2% |
| left_hip_roll_joint | 6.398e-01 | 2.510e-02 | 6.649e-01 | 3.8% |
| left_hip_yaw_joint | 1.296e-01 | 1.018e-02 | 1.398e-01 | 7.3% |
| left_knee_joint | 1.115e-01 | 2.510e-02 | 1.366e-01 | 18.4% |
| left_ankle_pitch_joint | 2.777e-03 | 7.219e-03 | 9.997e-03 | 72.2% |
| left_ankle_roll_joint | 3.871e-04 | 7.219e-03 | 7.607e-03 | 94.9% |

**Table S5**: **Motion Tracking Hyperparameters**

| Hyperparameter | Value |
|---|---|
| Architecture | |
| Actor MLP hidden dimensions | [512, 256, 128] |
| Critic MLP hidden dimensions | [512, 256, 128] |
| Activation function | ELU |
| Training | |
| Steps per environment | 24 |
| Max iterations | 30,000 |
| Learning rate | $1 \times 10^{-3}$ |
| Clip parameter | 0.2 |
| Entropy coefficient | 0.005 |
| Value loss coefficient | 1.0 |
| Discount factor ($\gamma$) | 0.99 |
| GAE $\lambda$ | 0.95 |
| Desired KL | 0.01 |
| Learning epochs | 5 |
| Mini-batches | 4 |

**Table S6**: **VAE Hyperparameters**

| Hyperparameter | Value |
|---|---|
| Architecture | |
| Latent dimension | 32 |
| Student encoder MLP hidden dimensions | [2048, 1024, 512] |
| Student decoder MLP hidden dimensions | [2048, 1024, 512] |
| Teacher hidden MLP hidden dimensions | [512, 256, 128] |
| Activation function | ELU |
| Training | |
| Learning rate | $5 \times 10^{-4}$ |
| Accumulated Gradient steps | 15 |
| KL loss coefficient | 0.01 |

**Table S7**: **Diffusion Policy Hyperparameters**

| Hyperparameter | Value |
|---|---|
| Architecture | |
| Horizon | 16 |
| Observation History | 4 |
| Embedding dimension | 512 |
| Attention heads | 8 |
| Transformer layers | 6 |
| Denoising steps | 20 |
| Training | |
| Batch size | 512 |
| Number of epochs | 1000 |
| Learning rate | $1 \times 10^{-4}$ |
| Weight decay | 0.001 |
| LR scheduler | Cosine |
| LR warmup gradient steps | 10,000 |
| EMA power | 0.75 |
| EMA max value | 0.9999 |

**Table S8**: **Motion Segments Tested in Sim and Real.**

| Name | Sim | Real [s] |
|---|---|---|
| Short Sequence | | |
| Cristiano Ronaldo (*33*) | Full | Full |
| Side Kick (*35*) | Full | Full |
| Single Leg Balance (*34*) | Full | Full |
| Swallow Balance (*34*) | Full | Full |
| Aerial Cartwheel | Full | Full |
| Double Kicks 1 | Full | Full |
| Double Kicks 2 | Full | Full |
| LAFAN1 (*83*) (about 3 minutes each) | | |
| walk1_subject1 | Full | [0.0, 33.0], [81.2, 86.7] |
| walk1_subject2 | Full | - |
| walk1_subject5 | Full | [146.7, 159.0], [206.7, 263.7] |
| walk2_subject1 | Full | - |
| walk2_subject3 | Full | [42.7, 75.7], [217.6, 230.6] |
| walk2_subject4 | Full | [154.4, 164.4], [218.6, 238.6] |
| dance1_subject1 | Full | [0.0, 118.0] |
| dance1_subject2 | Full | Full |
| dance1_subject3 | Full | - |
| dance2_subject1 | Full | - |
| dance2_subject2 | Full | - |
| dance2_subject3 | Full | [43.1, 163.1], [164.3, 184.3] |
| dance2_subject4 | Full | [156.3, end] |
| dance2_subject4 | Full | - |
| fallAndGetUp1_subject4 | Full | - |
| fallAndGetUp2_subject2 | Full | [0.0, 21.0], [74.0, 91.2], [94.0, 109.0] |
| fallAndGetUp2_subject3 | Full | [26.5, 46.5] |
| run1_subject2 | Full | [0.0, 50.0] |
| run1_subject4 | Full | - |
| run1_subject5 | Full | - |
| run2_subject1 | Full | [0.0, 11.0], [167.4, 204.4] |
| jumps1_subject1 | Full | [24.3, 42.3], [71.6, 81.6], [205.5, 226.5] |
| jumps1_subject2 | Full | - |
| jumps1_subject5 | Full | - |
| fightAndSports1_subject1 | Full | [16.8, 25.4], [201.6, end] |
| fightAndSports1_subject4 | Full | - |
| fight1_subject2 | Full | - |
| fight1_subject3 | Full | - |
| fight1_subject5 | Full | - |