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

Qiayuan Liao*[1], Takara E. Truong*[2], Xiaoyu Huang*[1], Guy Tevet[2], Koushil Sreenath†[1] and C. Karen Liu†[2]

*Abstract*— Learning skills from human motions offers a promising path toward generalizable policies for versatile humanoid whole-body control, yet two key cornerstones are missing: (1) a high-quality motion tracking framework that faithfully transforms large-scale kinematic references into robust and extremely dynamic motions on real hardware, and (2) a distillation approach that can effectively learn these motion primitives and compose them to solve downstream tasks. We address these gaps with BeyondMimic, a real-world framework to learn from human motions for versatile and naturalistic humanoid control via guided diffusion. Our framework provides a motion tracking pipeline capable of challenging skills such as jumping spins, sprinting, and cartwheels with state-of-the-art motion quality. Moving beyond simply mimicking existing motions, we further introduce a unified diffusion policy that enables zero-shot task-specific control at test time using simple cost functions. Deployed on hardware, BeyondMimic performs diverse tasks at test time, including waypoint navigation, joystick teleoperation, and obstacle avoidance, bridging sim-to-real motion tracking and flexible synthesis of human motion primitives for whole-body control. **https://beyondmimic. github.io/**.

Fig. 1. Overview of BeyondMimic: We train robust motion tracking policies $\pi$ for effective sim-to-real transfer, then perform offline distillation to learn a state-action diffusion model, and finally deploy the distilled policy on downstream tasks using guidance including waypoint navigation with obstacle avoidance (left) and joystick-based velocity control (right).

## I. INTRODUCTION

Learning from the rich diversity of human motion offers a promising path toward agile, human-like whole-body control in humanoid robots. Motion capture datasets contain diverse, dynamic behaviors, from stylish locomotion to complex whole-body motions, that, if transferred to robots, could form reusable building blocks for versatile, generalizable control across downstream tasks.

Recent advances in physics-based character animation [1], [2], [3], [4], [5] have already shown incredible progress to synthesize human motions into dynamic behaviors for whole-body control to solve downstream tasks. However, these successes remain limited to simulation, where agents enjoy idealized dynamics, unlimited actuation, and perfect observations, unlike real-world humanoids that face unmodeled dynamics, hardware limits, and imperfect state estimation. Achieving similar capabilities on hardware requires two missing capabilities: (1) a scalable, high-quality motion tracking framework that transforms kinematic references into robust, highly-dynamic motions while avoiding over-randomization, jitter, and degraded motion quality seen in
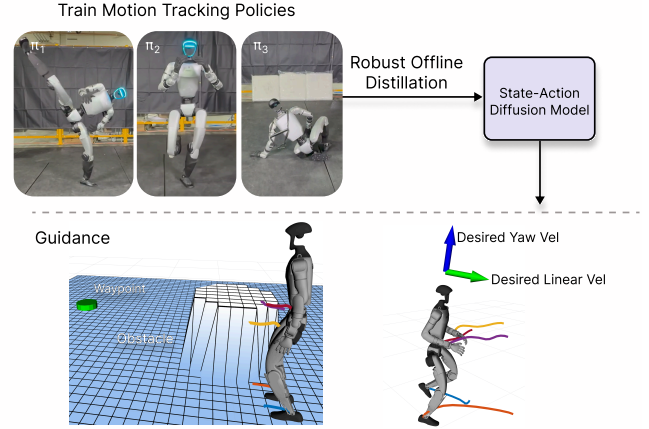
prior work, and (2) an effective sim-to-real recipe to distill learned motion primitives into a single policy capable of composing skills at test time for flexible, goal-driven control without retraining.

We introduce BeyondMimic, a real-world framework to address the above two challenges altogether. BeyondMimic first introduces a motion tracking pipeline that is able to perform highly dynamic motions, such as jumping spins, sprinting, and cartwheels, on real humanoid hardware with state-of-the-art motion quality. Beyond simply mimicking these human motions on humanoids, our framework further presents a unified guided-diffusion policy that synthesizes these motions or skills into novel motions for zero-shot, task-specific control at test time using only simple cost functions, solving downstream tasks flexibly without additional retraining.

We demonstrate BeyondMimic on a full sim-to-real pipeline: training robust tracking policies, distilling them into a diffusion-based controller, and deploying the result on physical hardware. Ultimately, our system performs a broad set of tasks, from waypoint navigation and joystick teleoperation to obstacle avoidance, all while preserving the naturalistic style and dynamic quality of the original human motions. By bridging gaps in sim-to-real transfer for both motion tracking and diffusion policy synthesis, BeyondMimic provides a practical foundation for advancing generalizable whole-body humanoid control. Our core

---

* Equal contribution; order decided by coin toss.

† Equal advising; order mirrors the coin toss, with the other lab listed last.

[1]Q. Liao, X. Huang, and K. Sreenath are with the University of California, Berkeley, CA 94720, USA. {qiayuanl, haytham.huang, koushils}@berkeley.edu.

[2]T.E. Truong, G. Tevet, and C. K. Liu are with Stanford University, Stanford, CA 94305, USA. {takarat, guytevet, ckliu}@stanford.edu.

contributions are:

- **Scalable Motion Tracking**: An open-source framework for training robust, high-quality, and highly dynamic motion tracking policies on real hardware, using a unified MDP and a single set of hyperparameters.
- **Guided Diffusion for Humanoids**: Demonstrating loss-guided diffusion for real-world humanoid whole-body control, enabling diverse downstream tasks via simple loss functions at test time, a capability not previously achievable with existing methods..
- **End-to-End Framework**: A complete pipeline from raw motion capture to hardware deployment, including motion tracking, diffusion distillation, and real-world deployment, as shown in Fig 1.

## II. RELATED WORK

### A. Motion Tracking

Early learning-based approaches to legged robot control focus on manually designing task-specific controllers [6], [7], [8], [9], which can produce robust locomotion but require extensive reward engineering for each task, lack naturalistic motion, and do not scale to the diverse skill set needed for general-purpose control. DeepMimic [10] offers an alternative by learning from human motion references to produce naturalistic, dynamically feasible behaviors while reducing the reward engineering burden, particularly suitable for humanoid robots. Since then, many tracking frameworks based on the DeepMimic paradigm have been developed.

Earlier motion-tracking approaches typically rely on a small set of similar reference motions to train single-task policies. These policies are co-trained with task-specific rewards alongside the DeepMimic objective, producing controllers capable of performing one task in one style. Examples include quadrupedal locomotion [11] and goalkeeping [12], and bipedal jumping and running [13]. To enable more dynamic skills, ASAP [14] proposes to use a real-to-sim pipeline that learns a delta action model from hardware experiments, improving simulation fidelity. However, this approach requires training motion-specific delta action models, which tend to overfit to the exact short motion. Similarly, KungfuBot [15] and HuB [16] achieve high-quality sim-to-real transfer with carefully hand-crafted domain randomization, but only for single short motions.

To move beyond single-motion policies, recent research explores scalable motion-tracking frameworks that can learn a diverse set of motions within a single policy. In physics-based animation, PHC [17] is a leading example of such a system, inspiring subsequent attempts in robotics to build general-purpose motion trackers. Early multi-motion trackers in robotics, such as OmniH2O [18], Exbody [19], and HumanPlus [20], demonstrate the feasibility of this approach but suffer significant motion quality degradation compared to their graphics counterparts. More recently, TWIST [21] enables good quality motion tracking but primarily for static motions, and similarly, CLONE [22] and UniTracker [23] focus mainly on low-dynamic walking motions with responsive upper-body tracking. On the other hand, GMT [24]

handles some dynamic motions but sacrifices global trajectory tracking in favor of relative velocity tracking and gait regularization reward to improve robustness. Grandia et al. [25] present perhaps the highest-quality multi-motion tracking to date, but only on a small robot and without dynamic motions.

To date, multi-motion tracking on real humanoid hardware with both high motion quality and support for highly dynamic skills has not been demonstrated. In this work, we aim to fill this gap by tracking diverse, minutes-long human references containing many distinct motions of varying style and difficulty.

### B. Diffusion in Robotics and Character Animation

Denoising diffusion models are emerging as powerful policy generators in robotics and animation, as they naturally handle multi-modal distributions and long-horizon sequences. In character animation, kinematic diffusion models offer intuitive conditioning through text [26], music [27], geometric constraints [28], or loss functions [29]. These planners are often paired with a separate low-level controller for physical execution [28], [30]. However, this two-stage approach often suffers from a "planning-control gap," where the planner produces out-of-distribution motions that the controller cannot robustly track [3]. This fragility is compounded by the challenge of selecting an appropriate replanning frequency: infrequent replanning struggles to adapt to environmental changes, while frequent replanning can prevent the agent from committing to a coherent strategy. While online replanning shows promise in simulation [30], it has yet to be validated in the real world.

To avoid the planning-control gap, another line of work focuses on learning a single, end-to-end policy that maps directly from state to an action distribution. This approach, Diffusion Policy [31], extends naturally to locomotion and character control. For example, DiffuseLoco [32], and BiRoDiff [33] learn unified policies that encode a wide range of skills, achieving smooth gait transitions and robust real-world deployment for quadruped robots. Meanwhile, PDP [34] extends the framework to physics-based animation, showing text prompt following and multi-modal behavior in push recovery. Using diffusion or flow matching policies in RL training also shows improved overall return and efficiency in motion tracking tasks [35]. While effective, such policies lack a mechanism for flexible test-time conditioning, such as loss-guided diffusion. This is because there is no straightforward way to compare a task goal defined in state space with an action sequence defined in joint space. As a result, conditioning on new tasks typically requires retraining.

A third category of methods seeks to combine the flexible guidance of planners with the robustness of end-to-end policies by diffusing a joint distribution over both states and actions. This enables test-time guidance through classifier-based rewards [36] or conditioning on future state goals [37]. Diffuser [36] pioneered this approach for offline reinforcement learning but demonstrates limited robustness. Decision Diffuser [37] similarly models state-action trajectories but

ultimately concludes that the actions lack robustness and replaces them with inverse dynamics post-processing. In contrast, Diffuse-CLoC [3] leverages robust offline distillation using PDP [34] and demonstrates strong performance in physics-based animation with this joint diffusion strategy. Nevertheless, applying guided, joint state-action diffusion models to complex real-world robotic control remains an open challenge, which we aim to solve in this work.

## III. SCALEABLE MOTION TRACKING

In this section, we detail the scalable pipeline for training motion-tracking policies. Using the same MDP and hyperparameters, the pipeline produces high-quality sim-to-real policies from minutes-long reference motions, demonstrating scalability across diverse motions. Note that since our tracking framework does not have a history, we omit timestep $t$ in the formulation for clarity. We use the subscript m to denote quantities from the reference motion. Unless otherwise stated, all amounts in this section are expressed in the world frame.

### A. Tracking Objective

We start from the retargeted reference motion, represented as keyframes $(\mathbf{q}_{\mathrm{m}}, \mathbf{v}_{\mathrm{m}})$ of generalized positions and velocities. Using forward kinematics, we obtain for each body $b \in \mathcal{B}$ its pose $T_{b,\mathrm{m}}$ and twist $\mathcal{V}_{b,\mathrm{m}}$, where $\mathcal{B}$ is the set of all robot bodies. The goal is to reproduce the reference motion on hardware with high fidelity in global coordinates.

Perturbations for robustness during training and the sim-to-real gap often lead to inevitable global drift. To preserve motion style while allowing such drift, the controller should not track absolute body poses.

Therefore, we select an anchor body $b_{\mathrm{anchor}} \in \mathcal{B}$, typically the root or torso, and anchor the motion reference to the desired tracking objective as follows: For the anchor body, use the reference motion directly: $\hat{T}_{b_{\mathrm{anchor}}} = T_{b_{\mathrm{anchor}},\mathrm{m}}$. For non-anchor bodies $b \in \mathcal{B} \setminus \{b_{\mathrm{anchor}}\}$, the desired pose is computed as $\hat{T}_b = T_\Delta T_{b_{\mathrm{anchor}},\mathrm{m}}^{-1} T_{b,\mathrm{m}}$, where $T_{b_{\mathrm{anchor}},\mathrm{m}}$ is the pose of the anchor body in reference motion, and $T_\Delta = (p_\Delta, R_\Delta)$ with $p_\Delta = [p_{b_{\mathrm{anchor}}.x}, \; p_{b_{\mathrm{anchor}}.y}, \; p_{b_{\mathrm{anchor}}.z,\mathrm{m}}]$ and $R_\Delta = R_z(\mathrm{yaw}(R_{b_{\mathrm{anchor}}} R_{b_{\mathrm{anchor}},\mathrm{m}}^\top))$. This hybrid transform shifts the motion into the robot's local frame by preserving height, aligning yaw, and translating the $xy$ origin under the robot. Finally, the desired twists remain unchanged, i.e., $\hat{\mathcal{V}}_b = \mathcal{V}_{b,\mathrm{m}}, \; \forall b \in \mathcal{B}$.

Robots often have many closely spaced bodies. Tracking all of them is therefore inefficient and often unnecessary. Instead, we select a subset of target bodies $\mathcal{B}_{\mathrm{target}} \subseteq \mathcal{B}$, forming the motion tracking objective as $g_{\mathrm{tracking}} = (\hat{T}_{b_{\mathrm{anchor}}}, \hat{T}_b, \hat{\mathcal{V}}_b), \; \forall b \in \mathcal{B}_{\mathrm{target}}$.

### B. Observations

We formulate the policy observation space as a single-timestep vector comprising three components. **(1) Reference phase.** We include the joint positions and velocities from the reference motion, $\mathbf{c} = [\mathbf{q}_{\mathrm{joint,m}}, \mathbf{v}_{\mathrm{joint,m}}]$, serving solely as phase information; the policy is not intended to track these

joint values directly. **(2) Anchor pose-tracking error.** We include the pose tracking error of the anchor body, $\xi_{b_{\mathrm{anchor}}} \in \mathbb{R}^9$, which consists of the three-dimensional position error and the first two columns of the rotation error matrix [38]. Since the reference motion is predefined in the world frame, this term implicitly provides orientation for balancing and global position for correcting drift. **(3) Other Proprioceptions.** We include the robot's root twist expressed in the root frame $^{b_{\mathrm{root}}}\mathcal{V}_{b_{\mathrm{root}}}$, joint positions $\mathbf{q}_{\mathrm{joint}}$ and velocities $\mathbf{v}_{\mathrm{joint}}$, and the previous action $\mathbf{a}_{\mathrm{last}}$.

The complete observation space is then $\mathbf{o} = [\mathbf{c}, \xi_{b_{\mathrm{anchor}}}, {}^{b_{\mathrm{root}}}\mathcal{V}_{b_{\mathrm{root}}}, \mathbf{q}_{\mathrm{joint}}, \mathbf{v}_{\mathrm{joint}}, \mathbf{a}_{\mathrm{last}}]$. Note that when position drift compensation is unnecessary or reliable state estimation is unavailable, the *linear* components may be omitted (i.e., the translational part of $\xi_{b_{\mathrm{anchor}}}$ and the linear component of $^{b_{\mathrm{root}}}\mathcal{V}_{b_{\mathrm{root}}}$).

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, $T_{b_{\mathrm{anchor}}}^{-1} T_b, \; \forall b \in \mathcal{B}$, enabling it to estimate body-wise tracking errors directly in Cartesian space.

### C. Joint Impedance and Actions

Adding joint impedance is a standard approach in animation and robotics. Many works in character animation [10], [17], [34] adopt high impedance for precise tracking, effectively reducing the control in free-space to a nearly kinematic problem. However, policies trained under such high-impedance settings are often impractical for hardware deployment, as they amplify sensor noise, reduce the passive compliance needed for impact absorption, and hinder implicit torque information from current and prior commands.

In comparison, we heuristically set joint stiffness and damping following Raibert et. al. [39]: $k_{\mathrm{p},j} = I_j \omega_n^2, \quad k_{\mathrm{d},j} = 2I_j \zeta \omega_n$, where $\omega_n$ is the natural frequency, $\zeta$ is the damping ratio, and $I_j = k_{\mathrm{g},j}^2 I_{\mathrm{motor},j}$ is the reflected inertia of the $j^{th}$ joint. We choose a damping ratio $\zeta = 2$ (overdamped) rather than 1 (critically damped) as in Raibert et. al. [39], 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 $10\,\mathrm{Hz}$, which promotes compliance with moderate gains.

The policy action is designed as normalized joint position setpoints: $\mathbf{q}_{j,t} = \bar{\mathbf{q}}_j + \alpha_j \mathbf{a}_{j,t}$, where $\mathbf{a}_{j,t}$ is the policy action output, $\bar{\mathbf{q}}_j$ is a constant nominal joint configuration, $\alpha_j = 0.25 \frac{\tau_{j,\mathrm{max}}}{k_{p,j}}$, with $\tau_{j,\mathrm{max}}$ representing the maximum allowable joint torque for joint $j$. This heuristic assumes that contacts generally occur around $\bar{\mathbf{q}}_j$ and that the robot hardware design ensures the maximum joint torque is proportional to the expected load. At low gains, these setpoints are not intended as desired position targets; rather, they serve as intermediate variables to generate desired torques, and are intentionally not clipped by joint kinematic limits.

### D. Rewards

We design the rewards in a simple, intuitive, and general way, consisting of (1) Task rewards as positive, uniformly

weighted, and expressed in task space; and (2) Minimal regularization penalties to avoid hurting tracking performance.

Task rewards are body tracking rewards. First, we compute error metrics for each target body $b \in \mathcal{B}_{\text{target}}$ based on the desired $(\hat{T}_b, \hat{\mathcal{V}}_b)$ and actual $(T_b, \mathcal{V}_b)$ poses and twists: $\mathbf{e}_{p,b} = \hat{\mathbf{p}}_b - \mathbf{p}_b$, $\mathbf{e}_{R,b} = \log(\hat{R}_b R_b^{\top})$, $\mathbf{e}_{v,b} = \hat{\mathbf{v}}_b - \mathbf{v}_b$, and $\mathbf{e}_{w,b} \approx \hat{\mathbf{w}}_b - \mathbf{w}_b$ assuming the orientaion error is small. Then, the mean squared errors are computed across all target bodies: $\bar{e}_\chi = \frac{1}{|\mathcal{B}_{\text{target}}|} \sum_{b \in \mathcal{B}_{\text{target}}} \|\mathbf{e}_{\chi,b}\|^2$, $\chi \in \{p, R, v, w\}$. Each error metric is then normalized using a Gaussian-style exponential function: $r(\bar{e}_\chi, \sigma) = \exp\left(-\bar{e}_\chi / \sigma_\chi^2\right)$, where $\sigma_\chi$ is a nominal error determined empirically.

The combined task reward is then defined as:

$$r_{\text{task}} = \sum_{\chi \in \{p,R,v,w\}} r(\bar{e}_\chi, \sigma_\chi). \tag{1}$$

For regularization, we include as few as three penalties critical for sim-to-real alignment. Among them, the joint limit penalty, $r_{\text{limit}}$, encourages joint positions to remain within the soft limits. The action rate penalty, $r_{\text{smooth}}$, encourages consecutive actions to be smooth, avoiding policies with excessive jitter. To penalize self-collision, we count the number of bodies where the self-contact force exceeds a predefined threshold as the total penalty, $r_{\text{contact}}$, over the bodies $b \notin \mathcal{B}_{\text{ee}} \subseteq \mathcal{B}$, where $\mathcal{B}_{\text{ee}}$ is the end-effector body set.

The total reward then becomes:

$$r = r_{\text{task}} - \lambda_l r_{\text{limit}} - \lambda_s r_{\text{smooth}} - \lambda_c r_{\text{contact}} \tag{2}$$

where $\lambda_l, \lambda_s, \lambda_c > 0$, define the reward weights.

Optionally, a global tracking reward for the anchor body $b_{\text{anchor}}$ can be added, following the same structure as $r_{\text{tracking}}$ but using errors $\mathbf{e}_{p,b_{\text{anchor}}}$ and $\mathbf{e}_{R,b_{\text{anchor}}}$.

*E. Termination and Reset*

An episode terminates under two conditions, indicating either a fall or tracking failure: (1) when the height or orientation (considering only pitch and roll) errors of the anchor body $b_{\text{anchor}}$ exceed predefined thresholds, or (2) when any height of the end-effector body $b \in \mathcal{B}_{\text{ee}}$ deviates significantly from the reference trajectory.

At each episode reset, the motion phase is adaptively sampled from the entire reference trajectory (see Sec. III-F). The robot is initialized at the corresponding reference configuration and velocity, with additional random perturbations to enhance robustness.

*F. Adaptive Sampling*

Training a long sequence of motion inevitably faces the problem that not all segments are equally difficult. Thus, uniform sampling over the entire trajectory, a common practice in prior works [14], [34], [17], often tends to oversample easy segments while undersampling hard ones, resulting in larger variance in rewards and training inefficiency.

Thus, it is natural to adaptively sample more frequently from harder regions. To achieve this, we divide the starting index of the entire motion into S bins of one second each and sample these bins according to empirical failure statistics.

Let $N_s$ and $F_s$ be the counts of episodes and failures starting in Bin s. The failure rate is smoothed over time using an exponential moving average to prevent discrete jumps in sampling caused by short-term fluctuations.

Since failures are more likely caused by suboptimal actions taken shortly before termination, we apply a non-causal convolution with an exponentially decaying kernel $k(u) = \gamma^u$, where $\gamma$ is the decay rate, to assign greater weight to recent past failures. The final sampling probability from Bin s is then:

$$p_s = \frac{\sum_{u=0}^{K-1} \alpha^u \bar{r}_{s+u}}{\sum_{j=1}^{S} \sum_{u=0}^{K-1} \alpha^u \bar{r}_{j+u}}, \tag{3}$$

where $\bar{r}_s$ is the smoothed failure rate of Bin s. We further mix the probability $p_s$ with a uniform distribution to preserve coverage of easier bins and mitigate catastrophic forgetting, using $p_s' = \lambda \frac{1}{B} + (1 - \lambda)p_s$, where $\lambda$ is the uniform sampling ratio. Starting bins are then drawn from $\text{Multinomial}(p_1', \ldots, p_S')$, prioritizing challenging regions.

*G. Domain Randomization*

We apply three domain randomization terms: the ground friction coefficient, default joint positions $\bar{\mathbf{q}}_j$ (for both actions and observations, effectively simulating joint offset calibration errors), and the torso's center of mass position. Additionally, we introduce perturbations in the environment during training to encourage the robot to learn policies robust to environmental variability.

IV. TRAJECTORY SYNTHESIS VIA GUIDED DIFFUSION

In this section, we present our state-action motion synthesis framework that enables task-specific control through guided diffusion. We first describe the training procedure, then detail the guidance mechanisms.

*A. Training*

We employ Diffuse-CLoC [3] to create a state-action co-diffusion framework capable of generating physically-grounded robot actions through guided sampling in state space. The model is trained to predict future trajectories, $\boldsymbol{\tau}_t = [\boldsymbol{a}_t, \boldsymbol{s}_{t+1}, \ldots, \boldsymbol{s}_{t+H}, \boldsymbol{a}_{t+H}]$, of a look-ahead horizon of $H$ timesteps, conditioned on an observation history $\boldsymbol{O}_t = [\boldsymbol{s}_{t-N}, \boldsymbol{a}_{t-N}, \cdots, \boldsymbol{s}_t]$ of $N$ steps of preceding state-action pairs.

The training itself follows a standard denoising diffusion process. First, in a forward pass, we incrementally add Gaussian noise to ground-truth trajectories. A denoising network, $x_{0,\theta}$, then learns to reverse this corruption. It takes a noised trajectory $\boldsymbol{\tau}_t^k$ as input and is trained to predict the original, clean trajectory $\hat{\boldsymbol{\tau}}_t^0 = x_{0,\theta}(\boldsymbol{\tau}_t^k, \boldsymbol{O}_t, \boldsymbol{k})$. The network's parameters are optimized by minimizing the Mean Squared Error (MSE) loss against the ground truth:

$$\mathcal{L} = \text{MSE}(x_{0,\theta}(\boldsymbol{\tau}_t^k, \boldsymbol{O}_t, \boldsymbol{k}), \boldsymbol{\tau}_t) \tag{4}$$

During inference, new trajectories are generated by starting with random noise and iteratively refining it using the

learned denoising function, following the DDPM update rule [40]:

$$\boldsymbol{\tau}_t^{\boldsymbol{k}-1} = \alpha_{\boldsymbol{k}}(\boldsymbol{\tau}_t^{\boldsymbol{k}} - \gamma_{\boldsymbol{k}}\epsilon_\theta(\boldsymbol{\tau}_t^{\boldsymbol{k}}, \boldsymbol{O}_t, \boldsymbol{k}) + \mathcal{N}(0, \sigma_{\boldsymbol{k}}^2 \boldsymbol{I})) \quad (5)$$

Here, $\epsilon_\theta(\cdot)$ is the predicted noise derived from the clean trajectory estimate $x_{0,\theta}(\cdot)$, $\alpha_{\boldsymbol{k}}$, $\gamma_{\boldsymbol{k}}$, and $\sigma_{\boldsymbol{k}}$ are DDPM coefficients. We use independent noise schedules for states and actions, denoted by $\boldsymbol{k} = (\boldsymbol{k_s}, \boldsymbol{k_a})$.

### B. Guidance

To direct the model toward specific objectives during inference, we use classifier guidance. Through Bayes' theorem [41], the score of a conditional distribution is decomposed into the model's learned unconditional score and a separate guidance term:

$$\nabla_{\boldsymbol{\tau}} \log p(\boldsymbol{\tau} \mid \boldsymbol{\tau}^*) = \nabla_{\boldsymbol{\tau}} \log p(\boldsymbol{\tau}) + \nabla_{\boldsymbol{\tau}} \log p(\boldsymbol{\tau}^* \mid \boldsymbol{\tau}), \quad (6)$$

where $\boldsymbol{\tau}^*$ is the optimal trajectory. In our application, this guidance term is a differentiable, task-specific cost function, $G_c(\boldsymbol{\tau})$. By establishing a relationship between the cost and the conditional likelihood, $p(\boldsymbol{\tau}^* \mid \boldsymbol{\tau}) \propto \exp(-G_c(\boldsymbol{\tau}))$, the guidance term simplifies to the negative gradient of the cost: $-\nabla_{\boldsymbol{\tau}} G_c(\boldsymbol{\tau})$.

This gradient provides a signal that drives the sampling process toward lower-cost regions of the trajectory space. This approach is highly versatile, enabling a single pre-trained model to be adapted to diverse objectives at inference time. We demonstrate this by evaluating our framework on three representative tasks, each defined by different types of trajectory constraints.

### C. Downstream Tasks

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:

$$G_{\mathrm{js}}(\boldsymbol{\tau}) = \frac{1}{2} \sum_{t'=t}^{t+H} \|V_{xy,t'}(\boldsymbol{\tau}_{t'}) - g_v\|^2, \quad (7)$$

where $\boldsymbol{\tau}$ is the predicted trajectory at timestep $t$, $V_{xy,t}$ pulls out the planar root velocity at timestep $t$ and $g_v \in \mathbb{R}^2$ is the goal root velocity from the joystick controller.

For the waypoint task, the cost function rewards proximity to the target while increasingly penalizing velocity as the agent gets closer, ensuring a stop upon arrival.

$$G_{\mathrm{wp}}(\boldsymbol{\tau}) = \sum_{t'=t}^{t+H} (1 - e^{-2d})\|P_x(\boldsymbol{s}_{t'}) - g_p\|^2 + e^{-2d}\|V_{x,t'}(\boldsymbol{\tau}_{t'})\|^2, \quad (8)$$

where $d = \|P_x(\boldsymbol{s}_{t'}) - g_p\|$ is the current scalar distance from the goal $g_p$, detached from the gradient computation.

For collision avoidance, we build a Signed Distance Field (SDF) to obtain the distance and gradient between body

positions and the nearest object, and define the cost function as:

$$G_{\mathrm{sdf}}(\boldsymbol{\tau}) = \sum_{t'=t}^{t+H} \sum_{b \in \mathcal{B}} B(\mathrm{SDF}(\mathbf{P}_{b,t'}(\boldsymbol{\tau})) - r_i, \delta), \quad (9)$$

where $\mathbf{P}_{b,t}(\boldsymbol{\tau})$ is the position of body $b$ at time $t$ given trajectory $\boldsymbol{\tau}$, $r_i$ is the approximation sphere radius of body $b$, and $B(x, \delta)$ is a relaxed barrier function [42]:

$$B(x, \delta) = \begin{cases} -\ln(x) & \text{if } x \geq \delta \\ -\ln(\delta) + \frac{1}{2}\left[\left(\frac{x-2\delta}{\delta}\right)^2 - 1\right] & \text{if } x < \delta \end{cases}. \quad (10)$$

## V. Motion Tracking Experiments and Results

This section presents the experimental validation of the motion tracking framework through both simulation and real-world deployment on the Unitree G1 humanoid robot, showing high-quality, robust and highly dynamic motions. Note that we train all policies with the same set of hyperparameters throughout this section.

### A. Sim-to-Sim Performance

We use the LAFAN1 dataset [43], retargeted by Unitree, which includes diverse and agile human motions such as sprinting, turning jumps, and crawling, along with short single-motion clips from prior works [16], [14]. The LAFAN1 dataset contains 40 several-minute-long references, each in a broad category with many distinct motions within that category. We randomly select 25 of these references, ensuring at least one from each category, and find that all can be completed in full during sim-to-sim evaluation. Due to the limited space of our test field, we purposely select 29 challenging clips from these references, each featuring dynamic and contact-rich behaviors, and test them on hardware to evaluate the sim-to-real transfer performance. The complete list of motions we have tested in sim-to-sim and sim-to-real transfer is available in Table I.

### B. Real-world Experiment Setup

We deploy our motion tracking policies, trained using the proposed method, on the Unitree G1 humanoid robot. All deployment code is written in C++ and optimized for real-time execution. Full-state estimation is provided at $500\,\mathrm{Hz}$ using a low-level generalized momentum observer combined with a Kalman filter [44]. 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) [45] for position correction or exclude state-estimation-dependent observations altogether.

All policies are executed onboard using ONNX Runtime [46] on the robot's CPU. Each inference step takes under $1.0\,\mathrm{ms}$, allowing seamless integration into the real-time estimation and control loop.

| Name | Sim | Real [s] |
|---|---|---|
| **Short Sequency** | | |
| Cristiano Ronaldo [14] | Full | Full |
| Side Kick [15] | Full | Full |
| Single Leg Balance [16] | Full | Full |
| Swallow Balance [16] | Full | Full |
| **LAFAN1 [43] (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 | - |

## C. Real-world Performance

We categorize and showcase four distinct classes of motion to demonstrate the superior performance, versatility, and generality of our motion tracking framework. Our evaluation includes challenging motions previously demonstrated by specialized frameworks, as well as novel sequences that have not, to our knowledge, been shown in prior works.

*a) Short Dynamic Motions Previously Demonstrated:* We begin with dynamic motions that have been featured in earlier works specialized in these motions, such as Cristiano Ronaldo's iconic celebration from ASAP [14] and a sidekick from KungfuBot [15]. Our system exhibits remarkable robustness in tracking these complex motions. Notably, while prior works demonstrate only single instances, our framework is capable of repeating these motions *five times in succession* by manually toggling the controller, without any degradation in stability or tracking quality.

*b) Static Motions Requiring Strong Balance:* We next consider motions that demand high levels of balance and postural control, including single-leg standing and the Swallow Exercise Balance from HuB [16]. Unlike HuB, which relies on task-specific domain randomization and parameter tuning for each motion, our framework completes these tasks using the *same hyperparameters* across all other motions. Although our policies may experience occasional recoveries, the motions can be completed successfully, highlighting the generality and robustness of our approach.

*c) Extremely Dynamic and Previously Undemonstrated Motions:* We push the limits of humanoid motion tracking by executing highly dynamic motions that have not, to our knowledge, been demonstrated by existing research (with the possible exception of Boston Dynamics Atlas). These include double- and single-leg hopping, two cartwheels in a row, out-and-back sprints, and forward jumps with $180°$ and $360°$ spins. Motions like the $360°$ spin jump are physically demanding even for adult humans, yet our policy performs them fluidly and reliably. Remarkably, these are not isolated clips, but parts of the *three-minute long references*, highlighting our framework's ability to scale to long motion sequences while preserving both precision and extreme agility hidden in the dataset.

*d) Stylized and Expressive Motions:* Last but not least, a fundamental objective of motion tracking is not only to execute motions stably, but to capture the stylistic elements that make them recognizably human. A framework should therefore be trained *not just to avoid falling*, but to preserve the timing, posture, and expression inherent in the original motion. Among the most stylistically rich categories in motion datasets are dancing and walking with diverse gaits. To this end, we demonstrate a wide variety of stylized motions, including the Charleston dance, Moonwalk, transitions from walking to crawling, elderly-style walking, and sport movements such as badminton and tennis. Our framework faithfully retains these stylistic features, producing motions that are *instantly recognizable* to casual observers. This highlights the capacity of our framework to go beyond stability to deliver high-fidelity, human-like behavior.

By successfully executing a wide spectrum of motions—ranging from static to highly dynamic, from well-studied to newly demonstrated, and from athletic feats to stylized motions—we show that our framework is robust and accurate, and *scalable, expressive, and broadly applicable*. These results demonstrate a unified motion-tracking system that attains this breadth and quality without task-specific tuning.

## D. Ablation on Adaptive Sampling

We ablate adaptive sampling during training to evaluate its effectiveness, using convergence speed as the metric. Convergence is defined as completing the entire motion in the sim-to-sim evaluation. As shown in Table II, adaptive sampling is critical for learning difficult segments in long motion sequences. Without it, all the long sequences listed failed sim-to-sim evaluation in mujoco [47] even after 30k training

| Motion | w/o AS | w/ AS |
|---|---|---|
| Christiano Ronaldo [14] | 3k | 1.5k |
| Swallow Balance [16] | 2.8k | 1.8k |
| dance1_subject1 | Failed (cartwheel) | 8k |
| dance2_subject1 | Failed (jump-spining) | 9k |
| fightAndSports1_subject1 | Failed (balance) | 10k |

TABLE III
SUCCESS RATE ON DIFFUSION STATE REPRESENTATION ABLATION

| State Representation | Walk + Perturb | Joystick Control |
|---|---|---|
| **Body-Pos State** | **100%** | **80%** |
| **Joint-Rot State** | 72% | 0% |

iterations, stalling at challenging motions such as balancing and cartwheels. With adaptive sampling, these sequences converged in only 10k iterations. Even for short motions like *Cristiano Ronaldo* from ASAP [14], adaptive sampling halved the required iterations (1.5k vs. 3k). These results demonstrate that adaptive sampling significantly accelerates training and improves robustness.

## VI. DIFFUSION EXPERIMENTS AND RESULTS

For the guided diffusion policy, we test key ablations for its sim-to-real transfer. In particular, we examine the role of the state representation and demonstrate that selecting the appropriate representation is critical for performance.

### A. Data

We use a subset of AMASS [48] and LAFAN1 [43] containing diverse walking motions. We train motion tracking controllers for each motion skill to generate action labels.

Since diffusion models require multiple denoising steps during inference, there is a significant enough delay between receiving an observation and commanding an action. To account for this latency, we include action delay domain randomization during training,

Once the trackers are trained, we follow a data collection procedure similar to PDP [34] and Diffuse-CLoC [3], where offline datasets are collected by rolling out the expert policy but with domain randomization added as a key difference.

### B. Experiment Setup

We train BeyondMimic using an observation history of $N = 4$ and predict a future horizon of $H = 16$. We limit the action prediction horizon to 8 steps through loss masking. The model uses a Transformer decoder with 6 layers, 4 attention heads, and a 512-dimensional embedding, totaling 19.95M parameters. Training is performed with 20 denoising steps and an attention dropout probability of 0.3.

At test time, the diffusion policy runs offboard using TensorRT on an NVIDIA RTX 4060 Mobile GPU. Inference is performed asynchronously in a separate thread due to latency; each step takes approximately 20 ms using 20 denoising steps. Gradients of the guiding cost are computed automatically using CppAD [49] within each denoising iteration. For the joint diffusion model, motion capture data is used to provide both environmental context for cost computation and improved state estimation in evaluation.

### C. Tasks and Metrics

We evaluate our method on two tasks:
- **Walk-Perturb:** During a 15-second walk, we apply instantaneous root velocity perturbations sampled uniformly from 0 to $0.5$ m/s once every second. We measure the fall rate in simulation.
- **Joystick Control:** We issue a sequence of directional commands: forward, backward, turn left, turn right—each lasting 3 seconds. Failure to follow commands or maintain balance is counted as a failure.

A fall is defined as any episode where the head height drops below $0.2$ m. Each experiment is run 50 times.

### D. State Representation Ablation

Choosing an appropriate state representation is critical for successful sim-to-real transfer. Local joint representations offer compactness and computational efficiency, making them attractive for real-time control applications. In contrast, predicting body positions directly requires larger representations and correspondingly larger models with increased inference time, but provides explicit spatial grounding in Cartesian coordinates. This tradeoff between computational efficiency and spatial interpretability raises the question of which encoding best supports diffusion learning while ensuring successful sim-to-real transfer. To address this question, we compare two representations:

- **Body Pos State**
  - *Global States:* Root position ($\mathbb{R}^3$), linear velocity ($\mathbb{R}^3$), and orientation represented as a rotation vector ($\mathbb{R}^3$). These are expressed relative to the character frame at the current timestep.
  - *Local States:* Cartesian positions ($\mathbb{R}^{3\mathcal{B}}$), and linear velocities ($\mathbb{R}^{3\mathcal{B}}$) of all bodies, expressed in the character frame at each timestep.
- **Joint Pos State**
  - *Global States:* Same as above—root position, velocity, and orientation as rotation vectors relative to the current frame.
  - *Local States:* Joint angles ($\mathbb{R}^J$), where $J$ is the number of joints, and joint velocities ($\mathbb{R}^J$).

Table III depicts the performance of each representation for each of the two evaluated tasks. We find that the Body-Pos state representation significantly outperforms Joint-Pos across both tasks. While Joint-Rot is theoretically more Markovian, small errors in joint position estimates accumulate through the kinematic chain, making this representation more susceptible to compounding diffusion prediction errors compared to directly predicting Cartesian body poses. Furthermore, adding guidance for joystick control on the

Joint-Rot state fails quickly and abruptly, even after careful tuning, likely because the base controller lacks robustness, and guidance pushes it further out of distribution.

## VII. DISCUSSION AND FUTURE WORK

### A. Sim-to-Real

Overall, we achieve strong sim-to-real performance across a wide range of motions for both tracking and the diffusion policy. The high-quality sim-to-real transfer likely stems from careful problem formulation and implementation, including accurate modeling with targeted randomization to avoid excessive domain randomization, design of observation and action space, and a low-latency C++ framework for communication and policy input/output processing.

By *carefully* addressing these key engineering details, we formulate the problem in a simple yet efficient way. These practices yield a small sim-to-real gap without requiring complex techniques such as real-to-sim adaptation or motion-specific randomization parameters, providing a strong foundation for future agile humanoid control. We will provide code for the motion tracking pipeline, where the reported performance can be reproduced on other robots of the same type.

Nonetheless, we still encounter occasional failures caused by state estimation drift, particularly when the end-effector contact assumption is violated, such as during the get-up motion for the tracking policy or in some diffusion policy's recovery mode. Developing a state estimator that generalizes across such diverse motions remains a challenging problem and is left for future work.

### B. Out-of-Distribution Behavior

A notable finding is that the diffusion model behaviour is inert compared to RL. In out-of-distribution scenarios, such as when the robot falls or is physically obstructed by a gantry, the robot tends to remain largely stationary. This allows human operators to simply and safely reposition the robot back up to be in distribution. This contrasts with typical RL policies, which often produce erratic or unstable behavior in such situations. The mild out of distribution behaviour exhibited by diffusion models could be highly beneficial for real-world deployment, where human safety and stability are critical.

### C. Skill Transitions

A key limitation of current action based diffusion models lies in their difficulty in transitioning between distinct skills. We can conceptualize these models as learning a manifold where each skill corresponds to a specific curve, with transitions occurring at the intersections between curves. While classifier guidance serves as a mechanism for directing the model to snap onto different curves, this approach becomes problematic when the target skill lies too far away on the manifold. In such cases, the model often becomes trapped in its current mode, unable to make the necessary transition. Future work could look into improving the ability for the model to transition between skills.

## REFERENCES

[1] Z. Luo, J. Cao, J. Merel, A. Winkler, J. Huang, K. Kitani, and W. Xu, "Universal humanoid motion representations for physics-based control," *arXiv preprint arXiv:2310.04582*, 2023.

[2] C. Tessler, Y. Guo, O. Nabati, G. Chechik, and X. B. Peng, "Masked-mimic: Unified physics-based character control through masked motion inpainting," *ACM Transactions on Graphics (TOG)*, vol. 43, no. 6, pp. 1–21, 2024.

[3] X. Huang, T. Truong, Y. Zhang, F. Yu, J. P. Sleiman, J. Hodgins, K. Sreenath, and F. Farshidian, "Diffuse-cloc: Guided diffusion for physics-based character look-ahead control," 2025. [Online]. Available: https://arxiv.org/abs/2503.11801

[4] L. Pan, Z. Yang, Z. Dou, W. Wang, B. Huang, B. Dai, T. Komura, and J. Wang, "Tokenhsi: Unified synthesis of physical human-scene interactions through task tokenization," in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 5379–5391.

[5] R. Yu, Y. Wang, Q. Zhao, H. W. Tsui, J. Wang, P. Tan, and Q. Chen, "Skillmimic-v2: Learning robust and generalizable interaction skills from sparse and noisy demonstrations," in *Proceedings of the Special Interest Group on Computer Graphics and Interactive Techniques Conference Conference Papers*, 2025, pp. 1–11.

[6] G. B. Margolis and P. Agrawal, "Walk these ways: Tuning robot control for generalization with multiplicity of behavior," in *Conference on Robot Learning*. PMLR, 2023, pp. 22–31.

[7] I. Radosavovic, T. Xiao, B. Zhang, T. Darrell, J. Malik, and K. Sreenath, "Real-world humanoid locomotion with reinforcement learning," *Science Robotics*, vol. 9, no. 89, p. eadi9579, 2024.

[8] Z. Xie, G. Berseth, P. Clary, J. Hurst, and M. Van de Panne, "Feedback control for cassie with deep reinforcement learning," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 1241–1246.

[9] N. Rudin, D. Hoeller, P. Reist, and M. Hutter, "Learning to walk in minutes using massively parallel deep reinforcement learning," in *Conference on robot learning*. PMLR, 2022, pp. 91–100.

[10] X. B. Peng, P. Abbeel, S. Levine, and M. van de Panne, "Deepmimic: Example-guided deep reinforcement learning of physics-based character skills," *ACM Trans. Graph.*, vol. 37, no. 4, pp. 143:1–143:14, Jul. 2018. [Online]. Available: http://doi.acm.org/10.1145/3197517.3201311

[11] X. B. Peng, E. Coumans, T. Zhang, T.-W. Lee, J. Tan, and S. Levine, "Learning agile robotic locomotion skills by imitating animals," *arXiv preprint arXiv:2004.00784*, 2020.

[12] X. Huang, Z. Li, Y. Xiang, Y. Ni, Y. Chi, Y. Li, L. Yang, X. B. Peng, and K. Sreenath, "Creating a dynamic quadrupedal robotic goalkeeper with reinforcement learning," in *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2023, pp. 2715–2722.

[13] Z. Li, X. B. Peng, P. Abbeel, S. Levine, G. Berseth, and K. Sreenath, "Reinforcement learning for versatile, dynamic, and robust bipedal locomotion control," *The International Journal of Robotics Research*, vol. 44, no. 5, pp. 840–888, 2025.

[14] T. He, J. Gao, W. Xiao, Y. Zhang, Z. Wang, J. Wang, Z. Luo, G. He, N. Sobanbab, C. Pan *et al.*, "Asap: Aligning simulation and real-world physics for learning agile humanoid whole-body skills," *arXiv preprint arXiv:2502.01143*, 2025.

[15] W. Xie, J. Han, J. Zheng, H. Li, X. Liu, J. Shi, W. Zhang, C. Bai, and X. Li, "Kungfubot: Physics-based humanoid whole-body control for learning highly-dynamic skills," *arXiv preprint arXiv:2506.12851*, 2025.

[16] T. Zhang, B. Zheng, R. Nai, Y. Hu, Y.-J. Wang, G. Chen, F. Lin, J. Li, C. Hong, K. Sreenath, and Y. Gao, "Hub: Learning extreme humanoid balance," *arXiv preprint arXiv:2505.07294*, 2025.

[17] Z. Luo, J. Cao, A. Winkler, K. Kitani, and W. Xu, "Perpetual humanoid control for real-time simulated avatars," in *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023, pp. 10 861–10 870.

[18] T. He, Z. Luo, X. He, W. Xiao, C. Zhang, W. Zhang, K. M. Kitani, C. Liu, and G. Shi, "Omnih2o: Universal and dexterous human-to-humanoid whole-body teleoperation and learning," in *Conference on Robot Learning*. PMLR, 2025, pp. 1516–1540.

[19] X. Cheng, Y. Ji, J. Chen, R. Yang, G. Yang, and X. Wang, "Expressive whole-body control for humanoid robots," *arXiv preprint arXiv:2402.16796*, 2024.

[20] Z. Fu, Q. Zhao, Q. Wu, G. Wetzstein, and C. Finn, "Humanplus: Humanoid shadowing and imitation from humans," *arXiv preprint arXiv:2406.10454*, 2024.

[21] Y. Ze, Z. Chen, J. P. AraÃšjo, Z.-a. Cao, X. B. Peng, J. Wu, and C. K. Liu, "Twist: Teleoperated whole-body imitation system," *arXiv preprint arXiv:2505.02833*, 2025.

[22] Y. Li, Y. Lin, J. Cui, T. Liu, W. Liang, Y. Zhu, and S. Huang, "Clone: Closed-loop whole-body humanoid teleoperation for long-horizon tasks," *arXiv preprint arXiv:2506.08931*, 2025.

[23] K. Yin, W. Zeng, K. Fan, Z. Wang, Q. Zhang, Z. Tian, J. Wang, J. Pang, and W. Zhang, "Unitracker: Learning universal whole-body motion tracker for humanoid robots," *arXiv preprint arXiv:2507.07356*, 2025.

[24] Z. Chen, M. Ji, X. Cheng, X. Peng, X. B. Peng, and X. Wang, "Gmt: General motion tracking for humanoid whole-body control," *arXiv preprint arXiv:2506.14770*, 2025.

[25] R. Grandia, E. Knoop, M. A. Hopkins, G. Wiedebach, J. Bishop, S. Pickles, D. Müller, and M. Bächer, "Design and control of a bipedal robotic character," *arXiv preprint arXiv:2501.05204*, 2025.

[26] G. Tevet, S. Raab, B. Gordon, Y. Shafir, D. Cohen-or, and A. H. Bermano, "Human motion diffusion model," in *The Eleventh International Conference on Learning Representations*, 2023.

[27] J. Tseng, R. Castellon, and K. Liu, "Edge: Editable dance generation from music," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023, pp. 448–458.

[28] Z. Xie, J. Tseng, S. Starke, M. van de Panne, and C. K. Liu, "Hierarchical planning and control for box loco-manipulation," pp. 1–18, 2023.

[29] K. Karunratanakul, K. Preechakul, S. Suwajanakorn, and S. Tang, " Guided Motion Diffusion for Controllable Human Motion Synthesis ," in *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*. Los Alamitos, CA, USA: IEEE Computer Society, Oct. 2023, pp. 2151–2162. [Online]. Available: https://doi.ieeecomputersociety.org/10.1109/ICCV51070.2023.00205

[30] G. Tevet, S. Raab, S. Cohan, D. Reda, Z. Luo, X. B. Peng, A. H. Bermano, and M. van de Panne, "Closd: Closing the loop between simulation and diffusion for multi-task character control," 2024. [Online]. Available: https://arxiv.org/abs/2410.03441

[31] C. Chi, Z. Xu, S. Feng, E. Cousineau, Y. Du, B. Burchfiel, R. Tedrake, and S. Song, "Diffusion policy: Visuomotor policy learning via action diffusion," *The International Journal of Robotics Research*, p. 02783649241273668, 2023.

[32] X. Huang, Y. Chi, R. Wang, Z. Li, X. B. Peng, S. Shao, B. Nikolic, and K. Sreenath, "Diffuseloco: Real-time legged locomotion control with diffusion from offline datasets," in *8th Annual Conference on Robot Learning*, 2024. [Online]. Available: https://openreview.net/forum?id=nVJm2RdPDu

[33] G. Mothish, M. Tayal, and S. Kolathaya, "Birodiff: Diffusion policies for bipedal robot locomotion on unseen terrains," 2024. [Online]. Available: https://arxiv.org/abs/2407.05424

[34] T. E. Truong, M. Piseno, Z. Xie, and C. K. Liu, "Pdp: Physics-based character animation via diffusion policy," in *SIGGRAPH Asia 2024 Conference Papers*, ser. SA '24. ACM, Dec. 2024, p. 1–10. [Online]. Available: http://dx.doi.org/10.1145/3680528.3687683

[35] D. McAllister, S. Ge, B. Yi, C. M. Kim, E. Weber, H. Choi, H. Feng, and A. Kanazawa, "Flow matching policy gradients," *arXiv preprint arXiv:2507.21053*, 2025.

[36] M. Janner, Y. Du, J. Tenenbaum, and S. Levine, "Planning with diffusion for flexible behavior synthesis," in *Proceedings of the 39th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato, Eds., vol. 162. PMLR, 17–23 Jul 2022, pp. 9902–9915. [Online]. Available: https://proceedings.mlr.press/v162/janner22a.html

[37] A. Ajay, Y. Du, A. Gupta, J. Tenenbaum, T. Jaakkola, and P. Agrawal, "Is conditional generative modeling all you need for decision-making?" 2023. [Online]. Available: https://openreview.net/forum?id=sP1fo2K9DFG

[38] Y. Zhou, C. Barnes, J. Lu, J. Yang, and 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.

[39] M. Raibert and F. Farshidian, "Workshop on whole-body control and bimanual manipulation: Applications in humanoids and beyond," Jun. 2025, presented at the Workshop on Whole-body Control and Bimanual Manipulation: Applications in Humanoids and Beyond, Robotics: Science and Systems (RSS) 2025.

[40] M. Welling and 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.

[41] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, "Score-based generative modeling through stochastic differential equations," in *International Conference on Learning Representations*, 2021.

[42] R. Grandia, F. Farshidian, R. Ranftl, and M. Hutter, "Feedback mpc for torque-controlled legged robots," 2019. [Online]. Available: https://arxiv.org/abs/1905.06144

[43] F. G. Harvey, M. Yurick, D. Nowrouzezahrai, and C. Pal, "Robust motion in-betweening," in *ACM Transactions on Graphics (Proceedings of ACM SIGGRAPH)*, vol. 39, no. 4. ACM, 2020.

[44] T. Flayols, A. Del Prete, P. Wensing, A. Mifsud, M. Benallegue, and O. Stasse, "Experimental evaluation of simple estimators for humanoid robots," in *2017 IEEE-RAS 17th International Conference on Humanoid Robotics (Humanoids)*. IEEE, 2017, pp. 889–895.

[45] K. Koide, M. Yokozuka, S. Oishi, and A. Banno, "Glim: 3d range-inertial localization and mapping with gpu-accelerated scan matching factors," *Robotics and Autonomous Systems*, vol. 179, p. 104750, 2024.

[46] O. R. developers, "Onnx runtime," https://onnxruntime.ai/, 2021, version: x.y.z.

[47] E. Todorov, T. Erez, and 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.

[48] N. Mahmood, N. Ghorbani, N. F. Troje, G. Pons-Moll, and M. J. Black, "AMASS: Archive of motion capture as surface shapes," in *International Conference on Computer Vision*, Oct. 2019, pp. 5442–5451.

[49] 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.