# Gait-Conditioned Reinforcement Learning with Multi-Phase Curriculum for Humanoid Locomotion

Tianhu Peng, Lingfan Bao and Chengxu Zhou[*]

*Abstract*—We present a unified gait-conditioned reinforcement learning framework that enables humanoid robots to perform standing, walking, running, and smooth transitions within a single recurrent policy. A compact reward routing mechanism dynamically activates gait-specific objectives based on a one-hot gait ID, mitigating reward interference and supporting stable multi-gait learning. Human-inspired reward terms promote biomechanically natural motions, such as straight-knee stance and coordinated arm-leg swing, without requiring motion capture data. A structured curriculum progressively introduces gait complexity and expands command space over multiple phases. In simulation, the policy successfully achieves robust standing, walking, running, and gait transitions. On the real Unitree G1 humanoid, we validate standing, walking, and walk-to-stand transitions, demonstrating stable and coordinated locomotion. This work provides a scalable, reference-free solution toward versatile and naturalistic humanoid control across diverse modes and environments.

## I. INTRODUCTION

Developing natural and efficient locomotion strategies for humanoid robots remains a core challenge in robotics. Unlike manipulators operating in structured environments, humanoids must adapt to dynamic settings where balance, adaptability, and smooth transitions are essential for real-world deployment. Human locomotion exhibits key biomechanical traits—straight-knee support, coordinated anti-phase arm swings, and smooth heel-to-toe contact—which collectively enhance balance, reduce energy expenditure, and regulate angular momentum [1]–[4].

These traits are not merely stylistic but serve fundamental functions. Rhythmic arm swing, for instance, helps cancel leg-induced angular momentum, improving trunk stability and reducing metabolic cost. While recent reinforcement learning (RL) work shows emergent anti-phase swing as a side effect of energy minimization [5], such methods rarely enforce angular momentum control directly. In contrast, we propose a biomechanics-inspired reward that explicitly penalizes residual angular momentum and promotes phase-symmetric arm-leg coordination—yielding stable and efficient gaits without relying on trajectory references.

Recent RL advances have enabled agile locomotion for legged robots [6]. However, reference-based methods such as Adversarial Motion Priors (AMP) [7] require large-scale motion capture datasets and often struggle with morphology
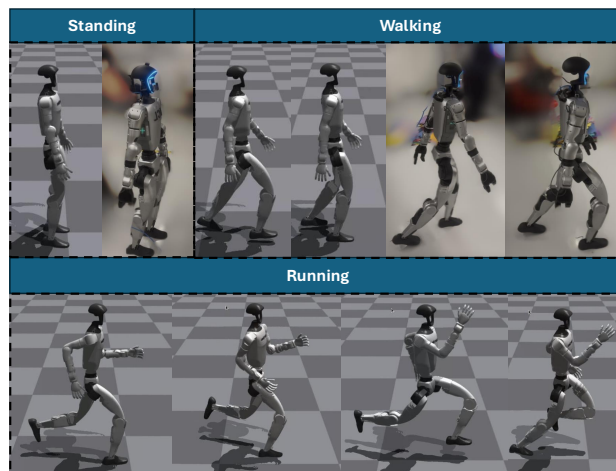


Fig. 1: Human-like multi-gait locomotion on the Unitree G1 humanoid, including standing, walking, and running. The learned reference-free policy exhibits straight-knee support, coordinated arm-leg motion, and natural transitions without MoCap references.

mismatch between human demonstrators and robot platforms. Moreover, their reliance on implicit imitation objectives limits interpretability, and prevents flexible reward design or task-specific modulation.

These limitations are amplified when integrating multiple gaits/skills into one controller. Prior approaches rely on multi-policy distillation [8], mixture-of-experts [9], or skill fusion [10], which typically require expert pretraining, switching mechanisms, and careful reward coordination across sub-policies. These architectures increase system complexity and training cost, and often suffer from interference between skill domains during deployment.

We propose a simpler alternative: a unified, gait-conditioned RL framework where a single recurrent policy learns standing, walking, running, and smooth transitions. A gait-conditioned reward routing mechanism activates gait-specific rewards based on a compact gait ID in the observation, mitigating interference and supporting stable multi-gait training.

To encourage human-like motion, we incorporate biomechanically grounded reward terms that promote straight-knee support, arm-leg coordination, minimal foot drag, and push-off dynamics. Our structured multi-phase curriculum enables progressive skill acquisition in simulation and robust deployment on hardware. This curriculum is inspired by bi-

The authors are with the Department of Computer Science, University College London, UK.

[*]Corresponding author, chengxu.zhou@ucl.ac.uk

ological motor development and bio-inspired gait adaptation studies [11]. A visual overview of the learned behaviors is shown in Fig. 1. In summary, our contributions are:

- A unified, reference-free RL framework for standing, walking, running, and transitions in a single recurrent policy.
- A gait-conditioned reward routing scheme that mitigates reward interference.
- Biomechanically grounded reward shaping for efficient, natural locomotion without MoCap.
- A progressive multi-phase curriculum for skill expansion and stable training.

## II. RELATED WORK

Learning natural and efficient locomotion has been a long-standing goal in RL for legged robots, especially in bipedal settings where stability and versatility are more demanding [6]. Existing methods can be broadly categorized into reference-based learning, reference-free approaches, modular multi-skill frameworks, and multi-behavior learning via reward or value function decoupling.

### A. Reference-Based Locomotion Learning

Early successes in humanoid locomotion were largely achieved through reference-based approaches, where policies imitate curated motion capture (MoCap) data. A seminal example is AMP [7], which introduced an adversarial framework to produce human-like motions without explicit trajectory tracking. In multi-gait settings, AMP uses several motion clips (e.g., trotting, pacing) with distinct root velocity and angular velocity profiles, switching between them based on commanded velocity. A velocity-tracking reward aligns the policy's motion with the reference clip's kinematics, effectively mapping different velocities to specific MoCap segments. AMP has been extended to quadrupeds [12], humanoid whole-body control [13], and cross-morphology transfer such as enabling quadrupeds to walk bipedally [14].

Beyond MoCap imitation, other reference-based methods leverage analytically generated trajectories or pre-defined gait libraries. Residual learning [15], [16] refines a nominal reference by predicting additive action offsets, enhancing adaptability while retaining its structure. Guided learning directly tracks analytical references such as Hybrid Zero Dynamics (HZD) [17] or Central Pattern Generator (CPG) templates [18], embedding rhythmic priors into training.

While effective, these methods depend on high-quality MoCap or analytical trajectories and often fail to generalize beyond their distribution. Morphological mismatches—differences in limb proportions, joint limits, or mass distribution—can degrade performance, and controllers trained on narrow gait libraries may struggle with unseen commands or disturbances.

### B. Reference-Free RL

Reference-free RL dispenses with predefined trajectories, instead optimizing handcrafted rewards to encourage energy-efficient, symmetric, and periodic gaits [19]. Early

work [20] demonstrated bipedal gait learning via symmetry and curriculum shaping; later studies targeted stepping stone navigation [21] and achieved sim-to-real transfer on Cassie with periodic rewards [16], extended to blind stair climbing [22] and vision-guided footstep placement [23], [24]. Heightmap-based perception further broadened terrain generalization [25].

While enabling greater generalization and design freedom, reference-free methods require careful reward tuning, converge slowly, and struggle with complex maneuvers such as agile transitions or jumping.

### C. Multi-Skill and Modular Architectures

Multi-skill RL aims to enable a single agent to execute diverse behaviors. Supervised strategies such as policy distillation [26] and DAgger [27] aggregate expert demonstrations to train a unified policy. For example, Han et al. [28] distill multiple task-specific controllers into a single terrain-adaptive locomotion policy, while Zhuang et al. [8] use DAgger to learn robust parkour behaviors from expert rollouts. These methods achieve strong performance but depend on extensive expert data and carefully managed curricula, limiting scalability to unseen behaviors or larger gait sets.

Modular approaches manage multiple skills via decomposed control structures. Mixture-of-Experts (MoE) controllers [9] dynamically select sub-policies, as in Perpetual Humanoid Control, which achieved real-time coordination of diverse human-like behaviors. Multiplicative Compositional Policies (MCP) [10] blend low-level primitives via multiplicative composition for flexible skill recombination. While modular frameworks offer high flexibility, they often incur significant architectural complexity and require skill-specific supervision or pretraining.

### D. Reward Routing and Value Function Decoupling

A key challenge in multi-behavior RL is mitigating reward interference across heterogeneous skills. Approaches such as MultiCritic Actor Learning [29] address this by maintaining a single actor with task-specific critic heads, improving stability. In locomotion, CPG- or phase-conditioned policies [18], [30] exploit rhythmic motion priors to stabilize multi-gait control.

While prior work like DeepMimic [31] incorporates one-hot skill identifiers into the observation to differentiate motions, these are used solely for policy conditioning without altering the reward structure. In contrast, our framework employs gait IDs both for conditioning and for dynamically routing gait-specific rewards, enabling structured multi-gait learning within a unified policy.

Despite these advances, achieving smooth and efficient gait transitions without explicit experts, motion references, or hierarchical planners remains an open challenge.

### E. Our Approach

Unlike AMP-based approaches [7], [31] that require large-scale motion capture datasets and often map commanded velocities to pre-recorded gait clips via velocity tracking

rewards, our method operates entirely without motion references. Instead of using MoCap-derived velocity profiles to drive gait switching, we define each gait—standing, walking, running, and transitions—directly through human-inspired biomechanical reward shaping, enabling the policy to discover natural transitions without clip-based supervision. In contrast to multi-policy or distillation frameworks [26], [27] that train and merge multiple specialized experts, we employ a single recurrent policy with a compact gait-ID encoding, allowing all gaits to be learned end-to-end within one unified architecture. This design eliminates expert-switching logic, reduces model complexity, and mitigates reward interference across gaits via a simple yet effective gait-conditioned reward routing mechanism. Compared to prior modular, distillation-based, or reference-driven methods, our framework emphasizes minimal architectural overhead, unified end-to-end training, and strong generalization across diverse locomotion behaviors, offering a scalable pathway toward robust and naturalistic humanoid locomotion control.

## III. METHODOLOGY

### A. Problem Setup

We address the problem of learning unified multi-gait control for humanoid locomotion, encompassing standing, walking, running, and transitions between these modes. The task is formulated as a partially observable Markov decision process (POMDP), defined by a tuple $(s_t, a_t, P, r_t, p_0, \gamma)$, where $s_t$ represents the latent system state at time $t$, $a_t$ is the action executed by the agent, $P$ denotes the transition dynamics, $r_t$ is the reward function, $p_0$ is the initial state distribution, and $\gamma$ is the discount factor. The objective is to learn a policy $\pi_\theta$ that maximizes the expected cumulative discounted return:

$$J(\theta) = \mathbb{E}_{\pi_\theta} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right] \tag{1}$$

### B. Policy Architecture

Our policy adopts an asymmetric actor–critic architecture [32], with both networks implemented as Long Short-Term Memory (LSTM) [33] layers followed by Multi-Layer Perceptrons (MLPs). The actor, designed for deployment, receives only proprioceptive observations, while the critic is given additional privileged information during training to improve value estimation.

The LSTM captures temporal dependencies across gait cycles, aiding balance recovery, handling delayed contact effects, and enabling smooth multi-gait transitions under partial observability. Actor inputs include base angular velocity and gravity orientation in the local frame, commanded velocities, joint position offsets from a nominal pose, joint velocities, the previous action, a sine–cosine gait phase encoding, and a one-hot gait ID (e.g., standing, walking, running, transition). The network outputs a 23-dimensional joint position offset vector, tracked by a joint-level proportional–derivative (PD) controller with fixed gains.

The critic augments the actor's inputs with privileged features such as foot contact states, ground reaction forces, contact friction/restitution, binary contact indicators at hip/knee joints, and externally applied disturbances with application points. This improves stability and accuracy in value estimation, even in noisy or partially observable settings.

An overview of the proposed gait-conditioned RL framework—highlighting the asymmetric architecture, gait-conditioned reward routing, and multi-phase curriculum—is shown in Fig. 2, illustrating how gait mode information conditions both actor and critic observations, from sensory input through network processing, reward computation, and final action selection.

### C. Gait-Conditioned Reward Routing

To enable multi-gait locomotion within a unified policy, we adopt a gait-conditioned reward routing mechanism. At each timestep, a gait mode ID is inferred based on the commanded velocity and the robot's dynamic state. This ID is encoded as a one-hot vector and appended to the policy's observation, conditioning the actor-critic network on the intended gait mode.

During training, a gait mask is applied to selectively activate mode-specific reward terms according to the current gait. As shown in Fig. 3, the complete reward vector includes shared components (e.g., task-level tracking and regulation rewards) as well as gait-related terms (e.g., contact patterns or push-off dynamics) specific to walking, running, or standing. The gait mask filters the reward vector such that only the relevant subset contributes to the learning signal. For example, running activates push-off and short contact duration rewards; walking enables swing height and foot symmetry terms; standing focuses on upright posture and stillness bonuses while disabling locomotion-specific terms.

This design ensures that each locomotion behavior—standing, walking, running, and their transitions—is optimized independently while coexisting within a single recurrent policy. It effectively mitigates reward conflicts, stabilizes multi-gait learning, and enables clean credit assignment across gait modes.

### D. Curriculum Learning Strategy

Inspired by the natural progression of human motor development—from static balance to walking and eventually running—we design a biologically grounded curriculum to enable stable acquisition of locomotion skills.

Training all gaits simultaneously often leads to reward conflict and unstable exploration. To address this, we adopt a multi-phase curriculum that gradually introduces gait modes, broadens the commanded velocity range, and incrementally activates coordination mechanisms.

The curriculum progresses along three dimensions:

- Gait complexity: Starting from walking, we progressively add standing, running, and transitional modes.
- Commanded velocity: The target velocity range expands from low-speed walking to high-speed running (up to 4.0 m/s).
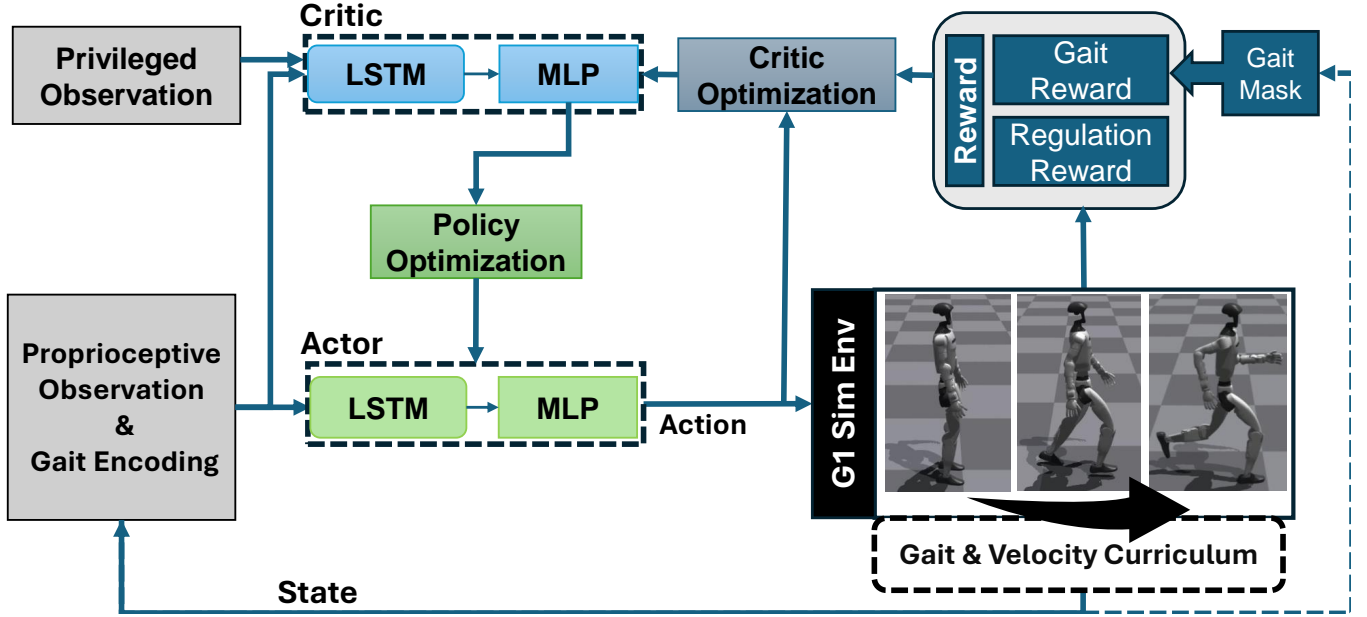
Fig. 2: Overview of our proposed gait-conditioned RL framework. A recurrent actor receives proprioceptive states with gait encoding. Gait-conditioned reward masks route mode-specific rewards, enabling multi-gait learning.



Fig. 3: Illustration of gait-conditioned reward routing during three example modes: Run, Walk, and Stand. The unified reward vector includes task-level and regulation rewards (shared across all gaits), as well as gait-specific terms (e.g., contact, push-off, or stability). A gait mask is used to activate only the relevant terms at each timestep. Ellipses (. . .) indicate omitted modes (e.g., transitions) or rewards.

- Coordination mechanisms: Reward routing, gait-aware masking, and smooth transition blending are gradually enabled.

We train the policy in three phases:

- **Phase 1** — Walking only: Only `Walk` (ID = 1) is enabled. The gait ID is fixed to `Walk` at all timesteps; other modes are disabled, and the reward mask activates walking-specific terms only. This stage develops periodic locomotion with stable contact, adequate foot clearance, and extended knees.
- **Phase 2** — Standing and Walk to Stand (W2S): We introduce gait-ID switching and corresponding rewards. When the commanded speed norm satisfies $\|\mathbf{v}_c\| < 0.1$ m/s, the system enters `W2S` (ID = 2). If low speed and double support persist for a manually specified 1.5 s, it switches to `Stand` (ID = 0). This hysteresis prevents premature switching.
- **Phase 3** — Running and Run-to-Walk (R2W): `Run` (ID = 3) is enabled using the Froude criterion $Fr = v^2/(g\,l)$ with $Fr > 0.5$. Upon speed reduction below this regime, the system enters `R2W` (ID = 4) until a manually specified 2.5 s stable duration completes, ensuring smooth velocity decay and contact reshaping. Running- and transition-specific rewards are activated, and all five modes (`Stand`, `Walk`, `W2S`, `Run`, `R2W`) are co-trained with full reward routing.

This progressive curriculum improves training stability, reduces reward interference, and facilitates generalization across diverse locomotion contexts.

### E. Human-Inspired Gait Design and Reward Shaping

While RL can produce stable locomotion behaviors, the resulting motions often appear overly crouched or energetically suboptimal. To encourage more natural and efficient gait patterns, we incorporate biomechanical principles observed in human locomotion directly into the reward design.

Specifically, our approach draws on the following insights:

- Straight-knee support improves force transmission and reduces muscular effort during stance phases [34].
- Anti-phase arm-leg coordination mitigates angular momentum buildup and stabilizes the upper body [1], [2].
- Gait transitions are gradual and biomechanically necessary—humans do not switch abruptly between running, walking, or standing, but instead adopt intermediate steps to reduce momentum and maintain balance. This motivates our inclusion of explicit transition gaits (e.g., walk-to-stand, run-to-walk), which align with observed human motor strategies [35], [36].

These principles are encoded through both gait-specific and shared reward components. Table I summarizes representative terms used to characterize and stabilize each locomotion mode. While the full reward set includes additional objectives (e.g., symmetry, torque minimization), we omit them here for brevity.

TABLE I: Representative Human-Inspired Reward Terms by Gait Mode

| Gait | Reward Term | Description |
|---|---|---|
| Walking | Contact Pattern | Encourage phase-aligned foot contacts based on cyclic gait timing. |
| | Foot Clearance | Promote sufficient foot lift during swing to avoid dragging. |
| | Straight Knee | Encourage extended knee during stance to improve support efficiency. |
| Running | Contact Pattern | Encourage alternating single-leg contact and flight phases. |
| | Push-Off Dynamics | Reward strong vertical and forward velocity during push-off. |
| | Short Contact | Penalize prolonged stance to promote dynamic running. |
| Standing | Contact Pattern | Encourage consistent double-foot support for static balance. |
| | Base Stability | Penalize base and joint motion to maintain upright posture. |
| Transition | Contact Pattern | Promote correct contact phasing during gait switching. |
| | Smooth Deceleration | Encourage gradual reduction in velocity to settle into stance. |

Key gait-specific rewards used to stabilize and differentiate locomotion behaviors. Additional terms such as symmetry and coordination are used but omitted for brevity.

*1) Angular Momentum Compensation via Arm-Leg Coordination:* Human arm swing plays a critical role in balancing the whole-body angular momentum during walking. Biomechanical studies [1]–[3] have shown that arm motion primarily serves to counteract leg-induced angular momentum—especially in the yaw direction—thereby enhancing trunk stability and reducing energy expenditure.

Inspired by this, we design a reward to promote human-like, dynamically balanced arm swing without relying on trajectory references. The total centroidal angular momentum $\mathbf{L}_{\text{total}}$ is computed via the standard decomposition [37]

To encourage coordinated arm-leg motion, we define the angular momentum reward as:

$$R_{\text{momentum}} = - \left( L_{\text{total},z}^2 \right) - 0.4 \left( L_{\text{la},z} - L_{\text{ra},z} \right)^2 \quad (2)$$

where $L_{\text{la}}, L_{\text{ra}} \in \mathbb{R}^3$ denote the angular momentum vectors of the left and right arms, and the subscripts $x$ and $z$ refer to pitch and yaw axes.

This reward encourages:

- *Whole-body momentum minimization:* Reduces residual angular momentum, particularly in the yaw direction;
- *Anti-phase yaw swing:* Promotes alternating arm motion to counteract leg-induced rotation.

### F. Training Platform

Training is performed in Isaac Gym [30], a GPU-accelerated physics simulator supporting parallel simulation of 1000 humanoid environments with NVIDIA RTX 4090

GPU. Policy optimization employs Proximal Policy Optimization (PPO) [38], with recurrent neural networks for both actor and critic, ensuring high sample efficiency and robust policy performance under partial observability.

### G. Domain Randomization

To improve policy robustness and facilitate sim-to-real transfer, we apply domain randomization during training. Several physical parameters are perturbed, including ground friction coefficient, robot base mass, and joint control properties. Additionally, external disturbances are introduced by applying randomized lateral pushes to the robot's base velocity at fixed intervals.

These perturbations expose the policy to a wide range of dynamics and environmental variations, encouraging generalizable locomotion behaviors that remain stable under uncertainty.

## IV. EXPERIMENT AND RESULTS

We evaluate our proposed gait-conditioned RL framework both in simulation (Isaac Gym) and real robot . The policy is trained to perform standing, walking, running, and smooth transitions between these modes using a single recurrent controller. Our experiments demonstrate that the learned policy achieves robust gait switching across varying commanded velocities, while exhibiting human-like motion and energy-efficient, coordinated behaviors—without relying on motion capture references or hierarchical planners.

### A. Simulation Results

*1) Gait Switching and Velocity Tracking:* To evaluate the controller's adaptability to varying locomotion demands, we command a velocity profile that ramps up and down between 0 and 3 m/s. Fig. 4 (bottom) shows that the actual forward velocity closely follows the commanded profile across a range of gaits. The top subplot visualizes contact phases along with background gait mode labels.

The policy exhibits stable transitions between standing (green), walking (blue), running (orange), and transitional gaits: W2S (purple) and R2W (red). These transitions occur smoothly without abrupt contact shifts, indicating robust internal gait modulation. In particular, the policy slows down naturally from running into walking , and settles into a balanced standing posture without external resets or mode switches.

*2) Angular Momentum Coordination:* Fig. 5 shows the Z-axis angular momentum of legs, arms, and the combined body over time. Throughout different gait modes—including transitions—the policy maintains low total angular momentum by coordinating the motion of arms and legs in anti-phase.

In walking and running, the arm and leg momenta are of opposite sign and similar magnitude, resulting in effective cancellation. During the R2W and W2S transitions, the coordination continues smoothly as the amplitude of each segment adapts to the underlying velocity and contact mode.
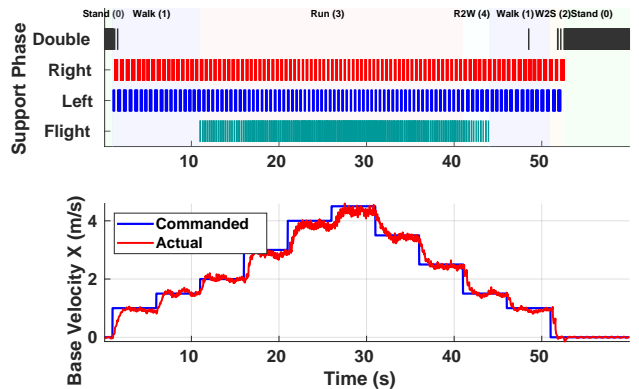


Fig. 4: Top: Support contact classification with background gait labels including Stand, Walk, Run, W2S, and R2W. Bottom: Commanded vs. actual forward base velocity. Gait transitions are smooth, dynamically consistent, and reflect internal coordination.
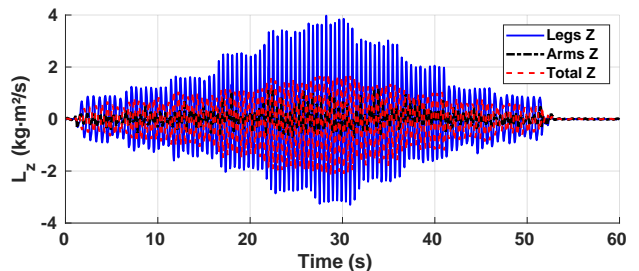


Fig. 5: Z-axis angular momentum of legs (blue solid), arms (red dashed), and total (black dash-dot). The policy achieves momentum coordination across standing, walking, running, and transitional gaits (W2S, R2W), with effective cancellation between arm and leg segments.

This demonstrates that the learned behavior is not only dynamically stable, but also biomechanically efficient.

These results confirm that our gait-conditioned policy learns to: (i) switch naturally between gaits; (ii) track commands robustly; and (iii) maintain biomechanical coordination.

*3) Ablations: Curriculum and Reward Routing:* We ablate two core components under identical architecture, hyper-parameters, and domain randomization: (i) *No Curriculum*: train all five gait modes (Standing, Walking, Running, W2S and R2W) jointly from scratch; (ii) *No Routing*: retain curriculum and gait ID in the observation, but disable reward routing so that all gait-specific terms (Table I) remain active simultaneously (mask removed). All settings are trained for an equal budget and evaluated on the same validation episodes, with results summarized in Table II.

Removing the curriculum forces the policy to learn all gait modes simultaneously from scratch, which greatly increases the optimization difficulty and prevents convergence within the same budget (episode length 31.26 vs. 972.6, return 2.45 vs. 89.3). Disabling reward routing, even with curriculum, activates mutually conflicting gait-specific rewards at every timestep (e.g., stillness from standing vs. forward velocity

TABLE II: Ablations under equal training budget. Len: mean episode length; Ret: mean return.

| Setting | Len. ↑ | Ret. ↑ |
|---|---|---|
| **Ours** (Curriculum + Routing) | 972.6 | 89.3 |
| No Curriculum (5 gait modes) | 31.26 | 2.45 |
| No Reward Routing (all terms active) | 11.59 | 0.00 |

from walking/running), causing unstable gradient updates and rapid performance collapse (episode length 11.59, return 0.00).

### B. Real-World Transfer

To evaluate the transferability of our policy from simulation to physical hardware, we deploy the learned controller on a real Unitree G1 humanoid robot. The deployed policy is directly transferred without any additional fine-tuning, sim-to-real adaptation, or dynamics randomization beyond what was already applied during training.

Our system successfully demonstrates stable standing, smooth walk-to-stand transitions, and walking on real hardware. These behaviors remain coherent and robust under moderate command perturbations and exhibit natural whole-body coordination. In particular, the arm-leg coordination and double-foot support in standing mode transfer well, indicating the effectiveness of our human-inspired reward shaping and gait-conditioned training framework.

Figure 6 shows snapshots of the deployed policy performing walking on the physical Unitree G1. The robot maintains biomechanically plausible postures such as extended knees and coordinated arm-leg motion, closely matching the simulated behavior.

Overall, our results validate that the proposed framework generalizes effectively across modalities and environments, with promising performance in real-world humanoid locomotion.

## V. CONCLUSION AND FUTURE WORK

We proposed a unified gait-conditioned RL framework for humanoid locomotion, combining gait-specific reward routing, biomechanically inspired reward shaping, and multi-phase curriculum learning. The resulting recurrent policy enables standing, walking, running, and smooth transitions, all within a single controller and without reliance on motion capture data.

Our method achieves robust and naturalistic multi-gait behaviors in physics-based simulation using a Unitree G1 humanoid. Initial real-world deployment demonstrates successful transfer of standing, walking, and walk-to-stand transitions. However, gait classification currently relies on manually defined gait IDs based on commanded velocity and Froude number, which constrains the emergence of novel gait patterns. The approach also depends on extensive manual design and tuning of numerous gait-specific reward terms, making it labor-intensive and increasingly challenging for larger or more complex gait sets. Moreover, as the number

of gaits grows, the shared policy may be prone to catastrophic forgetting or behavioral drift due to overlapping objectives.

We are currently working to extend this to dynamic running behaviors, which are already stable in simulation. Future directions include scaling to more complex scenarios such as traversal of uneven or deformable terrain, vision-guided locomotion in dynamic environments, and the incorporation of task-conditioned objectives for whole-body behaviors like climbing, object interaction, and human-robot collaboration.

### REFERENCES

[1] S. H. Collins, P. G. Adamczyk, and A. D. Kuo, "Dynamic arm swinging in human walking," *Proc. Royal Society B: Biological Sciences*, vol. 276, no. 1673, pp. 3679–3688, 2009.

[2] H. Pontzer, J. H. Holloway *et al.*, "Control and function of arm swing in human walking and running," *Journal of Experimental Biology*, vol. 212, no. 4, pp. 523–534, 2009.

[3] H. Herr and M. Popovic, "The roles of arm swing in human walking," *Journal of Experimental Biology*, vol. 211, no. 4, pp. 633–640, 2008.

[4] B. R. Umberger, "Effects of suppressing arm swing on kinematics, kinetics, and energetics of human walking," *Journal of Biomechanics*, vol. 41, no. 11, pp. 2575–2580, 2008.

[5] I. Radosavovic, A. Desai *et al.*, "Real-world humanoid locomotion with reinforcement learning," *Science Robotics*, vol. 9, no. 87, p. eadi9579, 2024.

[6] L. Bao, J. Humphreys *et al.*, "Deep reinforcement learning for bipedal locomotion: A brief survey," *arXiv preprint arXiv:2404.17070*, 2025.

[7] X. B. Peng, Z. Ma *et al.*, "Amp: Adversarial motion priors for stylized physics-based character control," *ACM Transactions on Graphics*, vol. 40, no. 4, pp. 1–20, 2021.

[8] Z. Zhuang, Z. Fu *et al.*, "Robot parkour learning," *arXiv preprint arXiv:2309.05665*, 2023.

[9] Z. Luo, J. Cao *et al.*, "Perpetual humanoid control for real-time simulated avatars," in *Proc. IEEE/CVF Int. Conf. on Computer Vision (ICCV)*, 2023, pp. 10 895–10 904.

[10] X. B. Peng, M. Chang *et al.*, "Mcp: Learning composable hierarchical control with multiplicative compositional policies," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 32, 2019.

[11] J. Humphreys and C. Zhou, "Learning to adapt through bio-inspired gait strategies for versatile quadruped locomotion," *arXiv preprint arXiv:2412.09440*, 2024.

[12] A. Escontrela, X. B. Peng *et al.*, "Adversarial motion priors make good substitutes for complex reward functions," in *Proc. IEEE/RSJ Int. Conf. on Intelligent Robots and Systems (IROS)*. IEEE, 2022, pp. 25–32.

[13] Q. Zhang, P. Cui *et al.*, "Whole-body humanoid robot locomotion with human reference," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2024, pp. 11 225–11 231.

[14] T. Peng, L. Bao *et al.*, "Learning bipedal walking on a quadruped robot via adversarial motion priors," in *Proc. Annual Conference Towards Autonomous Robotic Systems (TAROS)*, 2024, pp. 118–129.

[15] Z. Xie, P. Clary *et al.*, "Learning locomotion skills for cassie: Iterative design and sim-to-real," in *Proc. Conf. on Robot Learning (CoRL)*, 2020, pp. 317–329.

[16] J. Siekmann, Y. Godse *et al.*, "Sim-to-real learning of all common bipedal gaits via periodic reward composition," in *Proc. IEEE Int. Conf. on Robotics and Automation (ICRA)*, 2021, pp. 7309–7315.

[17] Z. Li, X. Cheng *et al.*, "Reinforcement learning for robust parameterized locomotion control of bipedal robots," in *Proc. IEEE Int. Conf. on Robotics and Automation (ICRA)*, 2021, pp. 2811–2817.

[18] J. Hwangbo, J. Lee *et al.*, "Learning agile and dynamic motor skills for legged robots," *Science Robotics*, vol. 4, no. 26, p. eaau5872, 2019.

[19] B. van Marum, M. Sabatelli, and H. Kasaei, "Learning perceptive bipedal locomotion over irregular terrain," *arXiv preprint arXiv:2304.07236*, 2023.

[20] W. Yu, G. Turk, and C. K. Liu, "Learning symmetric and low-energy locomotion," *ACM Transactions on Graphics*, vol. 37, pp. 1–12, 2018.

[21] Z. Xie, H. Ling *et al.*, "ALLSTEPS: Curriculum-driven learning of stepping stone skills," *Computer Graphics Forum*, vol. 39, pp. 213–224, 2020.

[22] J. Siekmann, K. Green *et al.*, "Blind bipedal stair traversal via sim-to-real reinforcement learning," in *Robotics: Science and Systems*, 2021.

Fig. 6: Real-world deployment of our learned policy on the Unitree G1 humanoid. The robot demonstrates smooth transitions between walking and standing, with natural limb coordination and upright posture—transferred directly from simulation without any fine-tuning.

[23] H. Duan, A. Malik *et al.*, "Sim-to-real learning of footstep-constrained bipedal dynamic walking," in *International Conference on Robotics and Automation*, 2022, pp. 10 428–10 434.

[24] ——, "Learning dynamic bipedal walking across stepping stones," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2022, pp. 6746–6752.

[25] B. Marum, M. Sabatelli, and H. Kasaei, "Learning vision-based bipedal locomotion for challenging terrain," *arXiv preprint arXiv:2309.14594*, 2023.

[26] A. A. Rusu, S. G. Colmenarejo *et al.*, "Policy distillation," *arXiv preprint arXiv:1511.06295*, 2015.

[27] S. Ross, G. Gordon, and J. A. Bagnell, "A reduction of imitation learning and structured prediction to no-regret online learning," in *Proc. Int. Conf. on Artificial Intelligence and Statistics (AISTATS)*, 2011, pp. 627–635.

[28] L. Han, Q. Zhu *et al.*, "Lifelike agility and play in quadrupedal robots using reinforcement learning and generative pre-trained models," *Nature Machine Intelligence*, vol. 6, no. 7, pp. 787–798, 2024.

[29] S. Mysore, G. Cheng *et al.*, "Multi-critic actor learning: Teaching rl policies to act with style," in *Proc. Int. Conf. on Learning Representations (ICLR)*, 2022.

[30] N. Rudin, D. Hoeller *et al.*, "Learning to walk in minutes using massively parallel deep reinforcement learning," in *Proc. Conf. on Robot Learning (CoRL)*, 2022, pp. 91–100.

[31] X. B. Peng, P. Abbeel *et al.*, "Deepmimic: Example-guided deep reinforcement learning of physics-based character skills," *ACM Transactions On Graphics (TOG)*, vol. 37, no. 4, pp. 1–14, 2018.

[32] M. Andrychowicz, B. Baker *et al.*, "Learning dexterous in-hand manipulation," in *The International Journal of Robotics Research*, vol. 39, no. 1. SAGE Publications, 2020, pp. 3–20.

[33] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[34] F. E. Zajac, R. R. Neptune, and S. A. Kautz, "Biomechanics and muscle coordination of human walking: Part i: Introduction to concepts, power transfer, dynamics and simulation," *Gait & Posture*, vol. 16, no. 3, pp. 215–232, 2002.

[35] A. E. Minetti and R. M. Alexander, "Translating resistive force theory to human gait transition: A non-linear optimization approach," *Journal of Experimental Biology*, vol. 203, pp. 2099–2108, 2000.

[36] T. Y. Hubel and J. R. Usherwood, "Transitions from walking to running: New insights from a single-body-center-of-mass model," *Journal of The Royal Society Interface*, vol. 10, no. 88, p. 20120977, 2013.

[37] J. Lee and A. Goswami, "Centroidal dynamics of a humanoid robot," *Autonomous Robots*, vol. 33, no. 3, pp. 291–311, 2012.

[38] J. Schulman, F. Wolski *et al.*, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.

# APPENDIX

## APPENDIX A: IMPLEMENTATION DETAILS

### A.1 Policy Architecture and Training Parameters

Our policy follows a recurrent actor-critic architecture with a single-layer LSTM module. The detailed configuration used for training is summarized in Table III.

TABLE III: Policy Architecture and PPO Training Hyperparameters

| Component | Setting |
|---|---|
| Policy architecture | LSTM + MLP |
| Actor hidden dimensions | [32] |
| Critic hidden dimensions | [32] |
| RNN type | LSTM |
| RNN hidden size | 64 |
| RNN layers | 1 |
| Activation function | ELU |
| Initial action noise std | 0.8 |
| Entropy coefficient | 0.01 |
| PPO iterations | 30,000 |

All training was conducted using the `rsl_rl` framework with a time step of 0.02s and a command update frequency of 10Hz.

### A.2 Reward Function Weights

We apply a gait-conditioned reward routing mechanism, where different subsets of reward terms are activated depending on the current gait mode. The primary reward weights used during training are listed in Table IV.

TABLE IV: Active reward terms and weights in the final multi-gait setting.

| Category | Reward Term | Weight |
|---|---|---|
| Regulation rewards | Vertical lin. vel. penalty ($v_z$) | $-2.0$ |
| | Horizontal ang. vel. penalty ($\omega_{xy}$) | $-0.05$ |
| | Orientation deviation | $-1.0$ |
| | Base height penalty | $-10.0$ |
| | DOF acceleration penalty | $-2.5 \times 10^{-7}$ |
| | DOF velocity penalty | $-1 \times 10^{-3}$ |
| | Collision penalty | $-10.0$ |
| | Action rate penalty | $-0.01$ |
| | DOF limit penalty | $-5.0$ |
| | Alive bonus | $0.15$ |
| | Hip pos. deviation | $-1.0$ |
| | Min. torso ang. vel. | $2.0$ |
| | Waist pitch deviation | $1.0$ |
| | Waist roll deviation | $1.0$ |
| | Waist yaw deviation | $1.2$ |
| | Torso yaw smoothness | $0.8$ |
| | Shoulder roll control | $3.0$ |
| Arm swing | Arm–leg momentum balance | $5.0$ |
| | Human-like arm swing energy | $0.3$ |
| | Elbow phase tracking | $2.5$ |
| | Arm swing symmetry | $2.0$ |
| | Arm swing–leg amp. match | $1.0$ |
| Walking | Feet swing height penalty | $-15.0$ |
| | Contact | $1.0$ |
| | Straight knee | $0.1$ |
| | Feet drag penalty | $-0.5$ |
| Standing | Contact (standing) | $2.5$ |
| | Base motion (standing) | $2.5$ |
| | Pose (soft upper) | $4.0$ |
| | Feet alignment | $0.5$ |
| | Uprightness | $1.0$ |
| | Feet drag penalty | $-0.2$ |
| | Feet flatness | $2.5$ |
| | Stillness bonus | $2.0$ |
| Walk→Stand (W2S) | Feet swing height penalty | $-20.0$ |
| | Contact | $1.0$ |
| | Smooth slowdown | $0.1$ |
| | Feet drag penalty | $-0.1$ |
| Running | Contact | $1.0$ |
| | Feet drag penalty | $-1.0$ |
| | Short ground contact | $0.2$ |
| | Feet swing height penalty | $-25.0$ |
| | Push-off reward | $1.0$ |
| | Forward velocity reward | $0.2$ |
| Run→Walk (R2W) | Contact | $1.0$ |
| | Smooth slowdown | $1.0$ |
| | Transition to walk speed | $0.5$ |
| | Feet swing height penalty | $-25.0$ |
| Task rewards | Tracking lin. vel. | $1.5$ |
| | Tracking ang. vel. | $1.0$ |