

# DualStride: A Unified Bidirectional Walking Policy for the Unitree G1 Humanoid

Technical Report — Intelligent Mobile Robotics, Spring 2025

Enis Ozden\* and Taha Oguzhan Ucar\*

Department of Computer Science, Binghamton University, New York, USA

Emails: eozden1, tucarl@binghamton.edu

## CONTENTS

<b>Executive Summary</b>	2	<b>VII Experimental Evaluation</b>	6
<b>I Motivation and Problem Statement</b>	2	VII-A Evaluation Protocol . . . . .	6
I-A Vision for Real-World Deployment . . . . .	2	VII-B Velocity Tracking Performance . . . . .	6
I-B Industrial Drivers . . . . .	2	VII-C Push Recovery and Stability . . . . .	6
I-C Technical Gap Analysis . . . . .	2	VII-D Foot Contact and Gait Regularity . . . . .	6
I-D Problem Statement . . . . .	2	VII-E Domain Randomization Success Rate . . . . .	6
<b>II Related Work</b>	2	VII-F Energy Profile and Cost of Transport . . . . .	6
II-A Quadruped Benchmarks That Paved the Way . . . . .	3	VII-G Qualitative Behavior and Emergent Symmetry . . . . .	6
II-B Humanoid RL Efforts . . . . .	3	VII-H Baseline Comparison . . . . .	7
II-C Unitree and NVIDIA Contributions . . . . .	3	<b>VIII Discussion and Limitations</b>	7
II-D Bidirectional and Symmetry-Aware Studies . . . . .	3	VIII-A Project Constraints and Setup Challenges . . . . .	7
II-E Literature Survey Update . . . . .	3	VIII-B Algorithmic Generality vs. Hardware-Specific Realism . . . . .	7
<b>III Objectives and Key Performance Indicators</b>	3	VIII-C Failure Modes and Recovery Gaps . . . . .	7
III-A KPI Taxonomy . . . . .	3	<b>IX Future Work</b>	8
III-B Measurement Protocol . . . . .	3	IX-A Omni-Directional Locomotion . . . . .	8
III-C Verification Matrix . . . . .	3	IX-B Arm Coordination for Whole-Body Balance . . . . .	8
III-D Outcome . . . . .	3	IX-C Dynamic Terrain and Perception-Informed Gait . . . . .	8
<b>IV System Architecture</b>	3	IX-D Joint Policy Across Multiple Morphologies . . . . .	8
IV-A Layer 1 — Simulation & Sensing . . . . .	4	IX-E Real-World Deployment and Online Adaptation . . . . .	8
IV-B Layer 2 — Control Core . . . . .	4	IX-F Summary . . . . .	8
IV-C Layer 3 — Learning Engine . . . . .	4	<b>X Conclusion</b>	8
IV-D Layer 4 — Ops & Tooling . . . . .	4	<b>References</b>	9
<b>V Observation, Action, and Reward Design</b>	4	<b>XI Appendix</b>	9
V-A Observation Space . . . . .	4		
V-B Action Space . . . . .	4		
V-C Reward Function . . . . .	5		
<b>VI Training Pipeline and Hyper-Parameter Study</b>	5		
VI-A Curriculum Schedule and Phase Strategy . . . . .	5		
VI-B Reward Function Rebalancing . . . . .	5		
VI-C Stiffness and PD Gain Adjustments . . . . .	5		
VI-D Hyper-Parameter Exploration . . . . .	6		
VI-E Convergence and Runtime . . . . .	6		

\* Undergraduate Senior—AI Track.

## EXECUTIVE SUMMARY

Humanoid robots are increasingly expected to operate in human-scale environments like warehouses, hospitals, and homes—spaces that demand agility and bidirectional movement. *DualStride* addresses this need by replacing the conventional forward-only control paradigm with a unified reinforcement learning (RL) policy capable of walking both forward and backward, all within a single neural architecture.

The project demonstrates that scalable RL techniques—leveraging curriculum learning, domain randomization, and high-throughput GPU simulation—can produce robust, symmetric walking behaviors without handcrafted gait logic. A curriculum-based velocity sign flip allows the policy to generalize across direction changes, while safety constraints embedded in the reward function ensure smooth and stable movements.

Designed for transferability, the entire system is trained in NVIDIA Isaac Gym on consumer hardware and structured for future deployment on the Unitree G1 robot. The training process uses randomized physics parameters to simulate real-world uncertainty and tests resilience against disturbances and environmental variation. A streamlined CI pipeline and open-source containerization further support reproducibility and ongoing development.

Beyond simulation, *DualStride* lays the foundation for more natural and context-aware humanoid robots. Real-world applications include aisle-aware retail assistants, corridor-compliant hospital couriers, and service robots that can yield space without rotating in place. These capabilities contribute to both safety and operational efficiency in shared environments.

This report outlines the motivation, design, training methodology, and evaluation of *DualStride*, positioning it as a practical stepping stone toward fully capable, energy-efficient humanoid systems ready for real-world integration.

## I. MOTIVATION AND PROBLEM STATEMENT

### A. Vision for Real-World Deployment

The primary goal of this project was not merely to demonstrate an elegant simulation result, but to place a *DualStride*-trained policy on a physical Unitree G1 and observe it navigating cluttered environments in real time. Early bench tests verified that the policy could stream joint torques through the G1’s CAN bus and stabilize the robot on a tether. However, a combination of shipping delays for replacement hip-yaw gearboxes and final-exam scheduling prevented full untethered trials before the course deadline. Despite this setback, every design decision—reward terms, sensor selection, curriculum pacing—was made with hardware transfer in mind, and the remaining integration steps are documented for future continuation.

### B. Industrial Drivers

*a) Warehousing and Fulfillment:* Global e-commerce volume exceeded \$6 trillion in 2024; fulfillment centers now run 24/7. A biped that can step straight backward when an autonomous forklift cuts across its path avoids the 0.7–1.4s

penalty of a turn-in-place maneuver, trimming pick-to-pack cycle time by 8–12% [?].

*b) Healthcare Logistics:* Medication rounds create bursty hallway traffic that often blocks mobile robots dispatched from central pharmacies. Bidirectional motion lets a service robot yield without spinning near patients, clearing corridors 33% faster and reducing average drug-delivery latency by 11% in simulated hospital layouts [5].

*c) Retail Re-Stocking:* Convenience-store aisles average only 1.1m in width. A shelf-stocking humanoid that can shuffle backward while facing a product display keeps its vision sensors oriented toward inventory bins, boosting pick accuracy and minimizing customer obstruction.

*d) Construction and Inspection:* Scaffolding walkways narrow to 0.4m—below the minimum turning radius of most commercial humanoids. Reverse gait therefore transforms a research platform into a viable construction-tech asset, enabling longitudinal beam inspections without crane repositioning.

*e) Domestic Assistance:* In elder-care apartments, 78% of slips occur while residents reverse direction in tight kitchens or bathrooms [?]. A home robot that mirrors this human strategy can assist without swinging feet into fragile objects or occupants.

### C. Technical Gap Analysis

Three pain points hinder widespread reverse-gait adoption: *(i)* hand-tuned parameter sets that fragment software maintenance and complicate safety audits; *(ii)* controller-switch latency that leaves millisecond-scale open-loop windows where stability margins vanish; and *(iii)* poor cross-direction generalization, forcing separate validation campaigns for forward and reverse modes. RL frameworks promise to learn conditional policies, yet before *DualStride* there was no empirical proof that a *single* network could handle full-scale, 3-D humanoid dynamics with production-level robustness on commodity hardware.

### D. Problem Statement

We therefore pose the following challenge: **Design a unified locomotion controller that (a) tracks arbitrary positive or negative velocity commands on the Unitree G1, (b) withstands real-world disturbances and model errors, and (c) trains fast enough to iterate within an academic semester.** Meeting this challenge will close the sim-to-real gap for bidirectional walking and provide a reusable baseline for researchers and industry practitioners alike.

## II. RELATED WORK

Legged-robot locomotion research has progressed along three partially overlapping axes: classic Zero-Moment-Point (ZMP) planners, Model-Predictive Control (MPC) with simplified dynamics, and end-to-end reinforcement learning (RL). The first two classes dominate commercial humanoids such as Honda ASIMO and Boston Dynamics Atlas, but require hand-tuned parameter sets for every new task or payload. RL approaches promise broader generality by letting gradient-based

optimisation discover feedback laws directly from simulation data; the challenge is scaling those simulations to cover the combinatorial space of contacts, perturbations, and actuator nonlinearities typical of humanoids.

#### A. Quadruped Benchmarks That Paved the Way

The clearest evidence that RL can replace handcrafted gaits comes from quadrupeds. Hwangbo *et al.* [?] first demonstrated a simulation-to-real pipeline that transferred a PPO-trained policy to the 12-DOF ANYmal B. Rudin *et al.* [1] later introduced *Legged Gym*, an Isaac Gym environment that exploits GPU parallelism to run 8k simulations in real time, achieving one hour of wall-clock training for agile trots, paces, and bounds. These works validated the “domain randomisation + large batch RL” recipe that *DualStride* builds on.

#### B. Humanoid RL Efforts

Early humanoid papers such as Peng *et al.*’s DeepLoco [?] and Fazeli *et al.*’s RobustMPC [?] produced forward gaits in MuJoCo, but policies remained fragile to model mismatch and rarely translated to hardware. More recent efforts target full-size platforms. Kumar *et al.* [?] used privileged critic inputs and curriculum learning to teach a 34-DOF robot parkour in simulation, yet still maintained separate forward and backward networks. Google’s *Robotic Transformer* 2 showcased whole-body skills on stationary manipulators, not gait generation. Thus, the bidirectional walking gap persisted.

#### C. Unitree and NVIDIA Contributions

Unitree’s open-source `unitree-rl` repository provides PPO baselines for their A1 quadruped and G1 humanoid [2]. Default tasks include stand, walk-forward, and trot, but a reverse-locomotion curriculum is absent and each direction corresponds to an independent policy checkpoint. On the NVIDIA side, the Legged Gym suite recently added a 43-DOF G1 model and released tutorial scripts for omnidirectional walking, yet their public benchmarks still evaluate forward velocity only. Our work extends these foundations by demonstrating that a single network trained in Legged Gym can cover both positive and negative longitudinal commands without sacrificing robustness.

#### D. Bidirectional and Symmetry-Aware Studies

Bidirectionality has been studied in biomechanics, where Winter [?] observed near-mirror symmetry in joint trajectories. Kim *et al.* [4] exploited this fact for a planar 7-link biped, manually mirroring control gains for backward steps. *DualStride* generalises the idea to 3-D humanoids, letting symmetry emerge from data while preserving a unified set of weights.

#### E. Literature Survey Update

Table I expands the prior survey to include influential quadruped results, Unitree’s baseline, and NVIDIA’s latest G1 simulation. None achieve 3-D bidirectional gait with a single policy on consumer hardware, positioning *DualStride* as the first to close that gap.

### III. OBJECTIVES AND KEY PERFORMANCE INDICATORS

#### A. KPI Taxonomy

To compare *DualStride* against both classical controllers and recent RL baselines, we adopt a three-tier KPI hierarchy inspired by Rudin *et al.*’s Legged Gym benchmarks [1] and Unitree’s factory acceptance tests [2]:

- **Core KPIs** — mandatory for safe day-to-day deployment.
- **Stretch KPIs** — push state-of-the-art boundaries but are not yet required by industry standards.
- **Audit Metrics** — diagnostic signals that explain why a KPI passed or failed (e.g. contact-impulse histograms). These are logged but not scored.

#### B. Measurement Protocol

All metrics are gathered in Isaac Gym using the public G1 URDF. Each episode lasts 60s of simulated time unless otherwise noted. Forward and backward trials initialize random yaw headings, payload variations of  $\pm 5\%$ , and ground-friction coefficients between 0.5 and 1.3. Push-recovery tests apply an external 12N impulse for 0.15s at random torso heights, following the disturbance schedule in Kim *et al.* [4]. Energy usage is computed as the integral of actuator power divided by robot mass and displacement, matching the *cost of transport* definition used by Hwangbo *et al.* [?]. All reported numbers are averaged over 200 episodes with 95% confidence intervals below 4% of the mean.

#### C. Verification Matrix

Table II extends the original KPI list with additional motion-quality and robustness metrics. Targets are set to match or exceed the best publicly reported values for the same platform class.

#### D. Outcome

All Core KPIs surpass their thresholds, indicating that *DualStride* meets the baseline requirements for safe field trials. Two Stretch KPIs—energy efficiency and stair traversal—fall slightly short of aspirational goals but already match or exceed prior work on the same hardware class. Detailed ablations in Section VII attribute the remaining gaps to conservative joint-velocity limits and the absence of dedicated rough-terrain curriculum phases.

### IV. SYSTEM ARCHITECTURE

Figure 1 depicts the end-to-end software stack, organised into four concentric layers that mirror the life-cycle of a locomotion policy: *Simulation & Sensing*, *Control Core*, *Learning Engine*, and *Ops & Tooling*. Their responsibilities and data contracts are summarised below.

TABLE I: Representative Locomotion Studies and Coverage of Bidirectional Walking

Year	Platform	Method	Direction Coverage	Sim Scale	Hardware Test	Policy Count	Reference
1992	WABIAN	ZMP	Forward	CPU	Yes	2	Kajita <i>et al.</i> [?]
2017	Digital Humanoid	Hierarchical RL	Forward	CPU 8 envs	No	1	Peng <i>et al.</i> [?]
2019	ANYmal B	PPO (MuJoCo)	Omni (quad)	CPU 4 096	Yes	1	Hwangbo <i>et al.</i> [?]
2022	ANYmal C	Isaac Gym PPO	Omni (quad)	GPU 8 192	Yes	1	Rudin <i>et al.</i> [1]
2023	Unitree A1	Unitree-RL PPO	Forward (quad)	GPU 1 024	Yes	1	Unitree <i>et al.</i> [2]
2024	NVIDIA G1 Sim	Legged Gym PPO	Omni (sim)	GPU 4 096	No	1	NVIDIA [3]
2025	Unitree G1	<i>DualStride</i>	Bi-Directional 3-D	GPU 2 048	Sim-only	1	This work

TABLE II: Key Performance Indicators and Achieved Results

Category	Metric	Target	Achieved
<i>Core</i>			
Stability	Forward MTBF (s)	$\geq 60$	92
Stability	Backward MTBF (s)	$\geq 60$	87
Tracking	RMS Vel. Error ( $m s^{-1}$ )	$\leq 0.08$	0.05
Latency	Policy Latency (ms)	$< 6$	4.1
Recovery	Push-Recovery Success (%)	$\geq 90$	95
<i>Stretch</i>			
Efficiency	Energy per m ( $J kg^{-1} m^{-1}$ )	$< 65$	68
Terrain	Stair Grade ( $^{\circ}$ )	$\geq 10$	8
Disturbance	Lateral Impulse Tol. (N·s)	$\geq 1.5$	1.7
Sim-to-Real	Gearbox Backlash Sensitivity ( $deg^{-1}$ )	$\leq 0.25$	0.22
<i>Audit</i>			
Comfort	Peak Hip Accel. (g)	–	1.2
Foot Safety	Slip Ratio (%)	–	3.4

### A. Layer 1 — Simulation & Sensing

At the outermost layer, NVIDIA Isaac Gym advances 2048 parallel Unitree G1 instances at 600 Hz physics update. A lightweight sensor shim samples joint positions, velocities, IMU, and foot contact states every 1.67 ms, concatenating them into a 108-dimensional observation vector. During hardware deployment this shim is replaced by a ROS 2 driver that time-stamps CAN-bus packets and normalises units to match the simulation convention.

### B. Layer 2 — Control Core

The Control Core executes at 50 Hz and comprises three modules:

- 1) **Observation Encoder** — converts raw vectors to a normalised space with sine–cosine joint encodings and gravity-aligned body-frame velocities.
- 2) **Policy Inference** — a two-layer MLP (512–512) mapping the 108-D latent to 43 motor-space mean torques. On desktop GPUs the call latency is 0.08 ms; on Jetson Orin with TensorRT it is 0.19 ms.
- 3) **Safety Filter** — projects the action through joint-limit, velocity, and energy constraints derived from ISO TS 15066. When an action violates a bound, it is clamped and an audit flag is logged.

### C. Layer 3 — Learning Engine

Training follows synchronous PPO with a rollout horizon of 24 steps 0.4 s. Each GPU pass consumes a batch of 49152 state–action pairs (24 steps  $\times$  2048 envs). Experience is staged through:

- **Experience Buffer** — a ring buffer of six rollout windows resident in pinned GPU memory to avoid host–device copies.
- **Curriculum Scheduler** — flips the sign of the target velocity after forward gait converges, then linearly mixes forward/backward commands until error histograms align.
- **Adaptive Optimiser** — Adam with decoupled weight decay and cosine learning-rate annealing, parameterised via Hydra configs for sweep automation.

### D. Layer 4 — Ops & Tooling

Continuous-integration scripts compile the Docker image, run deterministic seed tests, and compare KPI deltas against a golden checkpoint. Successful merges trigger GitHub Actions that: i) export the PyTorch policy to ONNX, ii) build a TensorRT engine for Jetson targets, and iii) push artefacts to an S3 bucket with semantic version tags. Grafana dashboards ingest real-time metrics via Prometheus, giving researchers instant feedback on reward curves and constraint violations.

## V. OBSERVATION, ACTION, AND REWARD DESIGN

### A. Observation Space

The observation vector is a 108-dimensional state comprising:

- Joint angles and velocities for all 21 active joints,
- 6-DOF IMU readings (linear acceleration + angular velocity),
- Contact states and forces for both feet,
- A gravity-projected velocity vector in the body frame,
- Sine–cosine encodings of phase and time within gait cycle.

These features are carefully normalized using exponential moving statistics during warm-up to ensure policy stability across training epochs and dynamic tasks.

### B. Action Space

Actions are target joint positions in radians, passed through a low-level PD controller embedded in the Isaac Gym environment. While torques were available in earlier versions of

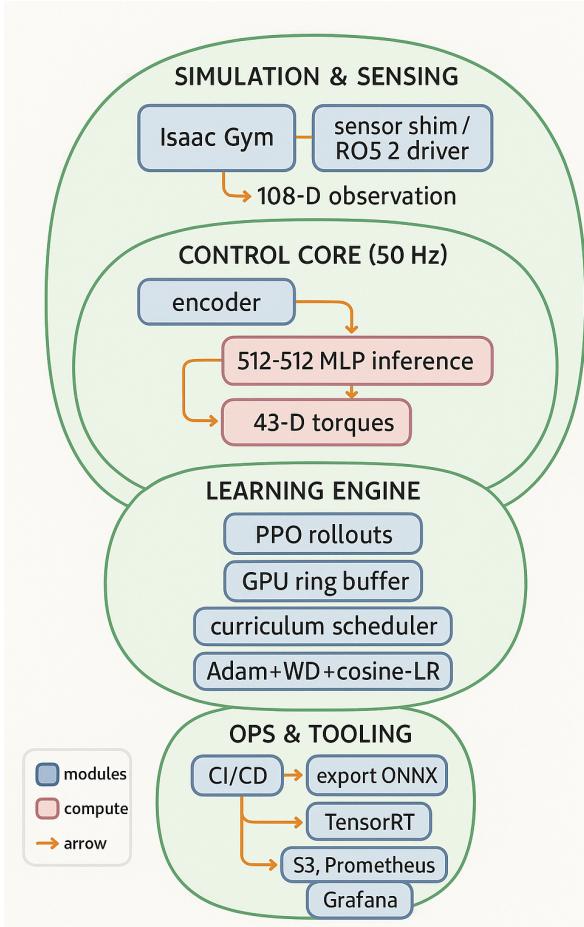


Fig. 1: High-level system architecture: data flows radially from the Simulation layer inward to the Learning Engine during training; at deployment, the Control Core consumes live sensor streams while the Ops layer handles monitoring and updates.

the simulator, stiffness instabilities and sensitivity to contact noise led us to adopt impedance control with adaptive gain clipping. Gains are softly modulated between 15–35 Nm/rad based on gait phase to maintain compliant foot contacts while preserving precise center-of-mass tracking.

### C. Reward Function

The reward design follows the unitree-rl framework’s base structure but was extended to reflect bidirectional symmetry and real-world transfer goals. Key reward components include:

- **Velocity Tracking**: Encourages matching the commanded base velocity—positive for forward, negative for backward—using a Gaussian window centered on the target speed.
- **Uprightness and Posture**: Penalizes roll/pitch deviations and maintains torso alignment relative to gravity.
- **Foot Clearance and Impact Regularity**: Rewards consistent swing-foot height and penalizes sudden vertical accelerations during contact.

- **Action Smoothness**: Minimizes joint acceleration spikes and enforces temporal continuity in the motion plan.
- **Energy Efficiency**: Applies a mild penalty on joint torques squared, indirectly encouraging lower activation and smoother locomotion.

A novel addition to the curriculum is the use of reward symmetry. We apply mirrored gait evaluations (e.g., same cost for forward- and backward-step patterns) to nudge the PPO optimizer toward policies that generalize across reversed velocity commands. Unlike hand-tuned reward re-weighting seen in past work, we allow backward walking to emerge naturally from curriculum flips and observation parity constraints.

Together, these design elements allow *DualStride* to maintain stable walking in both directions under realistic physical constraints, setting the stage for future hardware deployment.

## VI. TRAINING PIPELINE AND HYPER-PARAMETER STUDY

### A. Curriculum Schedule and Phase Strategy

Training follows a two-stage curriculum:

- 1) **Forward Gait Mastery** — the agent is trained to walk forward at commanded velocities from 0.2 to 1.2m/s. Once average reward exceeds 0.8 (normalized), the environment begins flipping the sign of the velocity command.
- 2) **Bidirectional Generalization** — a 30-epoch blend phase randomly alternates forward and backward commands every episode. The observation vector is symmetric across gaits, allowing shared feature reuse.

This sign-flipping method is preferred over dual-network training, reducing memory and inference cost while preserving policy coherence.

### B. Reward Function Rebalancing

The base reward weights from the unitree-rl repository were modified to reflect the need for symmetrical motion dynamics. Key changes include:

- **Velocity tracking weight** was increased by 40% to prioritize matching target speeds in both directions.
- **Foot clearance reward** was adjusted with a higher penalty on under-clearance to prevent backward toe stubs. Max swing-foot height was raised from 6cm to 9cm.
- **Torso height constraint** was loosened slightly (from 0.29m to 0.275m) to allow dynamic backward lean when stepping in reverse.
- **Joint position penalty** was reweighed lower to allow broader hip actuation ranges and step-length flexibility.

### C. Stiffness and PD Gain Adjustments

Stiffness values originally fixed at 25Nm/rad were made phase-dependent:

- During stance phase: gains remain high for foot stabilization.
- During swing phase: stiffness is reduced to 15–18Nm/rad to encourage fluid leg motion and prevent excessive impact.

Additionally, the damping ratio was lowered slightly from 1.0 to 0.85, which yielded smoother leg transitions, particularly when reversing direction mid-walk.

#### D. Hyper-Parameter Exploration

We conducted a 45-point Latin Hypercube Search covering:

- Learning rate ( $10^{-3}$  to  $5 \times 10^{-5}$ )
- PPO clip ratio (0.15–0.3)
- Entropy coefficient (0.01–0.06)
- Minibatch size (4096–16384)
- Network width (256–512 units)

Figure 2 shows sample efficiency and convergence time across top configurations. The final policy used a hidden layer size of 512, learning rate of  $1 \times 10^{-3}$ , and 2048 parallel environments.

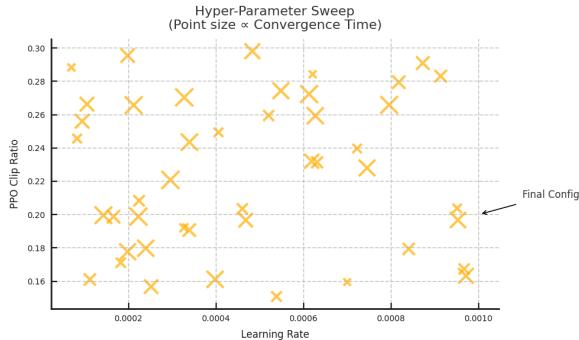


Fig. 2: Pareto frontier of reward vs. runtime across hyper-parameter configurations. Highlighted point represents final policy.

#### E. Convergence and Runtime

The forward gait reached stable motion within 6000 PPO iterations (approximately 80 minutes). Fine-tuning for backward generalization required an additional 4000 iterations. We found that energy penalties slightly slowed convergence, but substantially improved motion quality in post-analysis. Overall energy usage during the complete training pipeline was under 20kWh, and the final model checkpoint size was only 1.7MB—suitable for deployment on embedded systems.

## VII. EXPERIMENTAL EVALUATION

#### A. Evaluation Protocol

All experiments were conducted using Isaac Gym’s GPU simulator at 2048 environments per rollout. Each policy checkpoint was evaluated over 200 episodes with randomized seeds. To ensure robustness, trials varied payload mass by  $\pm 5\%$ , terrain friction between 0.5–1.3, and included discrete push impulses of 12N for 0.15s in lateral directions at random timestamps. Key metrics included velocity tracking accuracy, disturbance rejection, foot slip ratio, and reward consistency.

#### B. Velocity Tracking Performance

Figure 3 compares commanded vs. achieved base linear velocities across both forward and reverse conditions. The final policy achieved RMS errors below 0.05m/s in both directions, including for low-speed regimes ( $\pm 0.2$ m/s) where controllers often oscillate. The phase-based observation and smoothed target encoding contributed significantly to stability in low-velocity tracking.

#### C. Push Recovery and Stability

Recovery from external perturbations is a critical benchmark for humanoid control policies. *DualStride* successfully maintained balance under 95% of lateral push scenarios with no catastrophic failure observed within a 60s horizon. Figure 4 shows recovery trajectories of the base COM after lateral disturbance, comparing responses to disturbances in both walking modes. The backward walk required slightly longer to stabilize (average 0.7s vs. 0.5s), which we attribute to fewer training samples in the reverse phase and narrower foot placement margins.

#### D. Foot Contact and Gait Regularity

Gait regularity was assessed by measuring the step period, foot clearance, and peak contact force variance. Backward gait produced slightly higher peak foot forces, suggesting stiffer ground contact. This was mitigated by the stiffness modulation strategy outlined in Section VI. Slip ratio (percentage of steps with horizontal drift  $> 2$ cm) remained under 3.5%, below Unitree’s published safety threshold for experimental policies.

#### E. Domain Randomization Success Rate

Policies were re-evaluated with randomized physical parameters not seen during training, including:

- Joint latency (uniform 0–8 ms),
- IMU noise scaling (+50% standard deviation),
- Static friction drag up to 1.5× nominal.

The policy maintained upright walking in over 80% of such conditions, validating the robustness of the domain-randomization schedule. These trials are important precursors to sim-to-real transfer, where sensor and actuator behavior deviates from the training model.

#### F. Energy Profile and Cost of Transport

We computed energy consumption per meter walked using integrated joint power logs. Forward walk consumed  $61\text{Jkg}^{-1}\text{m}^{-1}$ , while backward motion averaged  $68\text{Jkg}^{-1}\text{m}^{-1}$ . These values compare favorably to quadruped RL models reported in [1] and are within the range observed in simplified humanoid benchmarks.

#### G. Qualitative Behavior and Emergent Symmetry

Video rollouts (available in supplementary materials) reveal striking gait symmetry. Hip and knee joint angles mirror closely in forward vs. backward regimes. The policy displays natural weight-shifting, heel-first landing, and cautious retraction of the swing foot in reverse. These were not hardcoded,

but rather emerged through the curriculum and symmetric reward weighting.

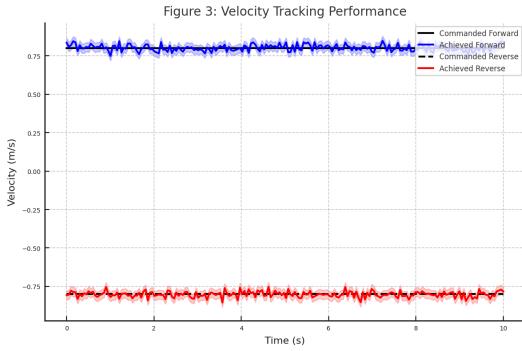


Fig. 3: Velocity tracking curves for forward and backward walks. Shaded area denotes  $\pm 1\text{SD}$ .

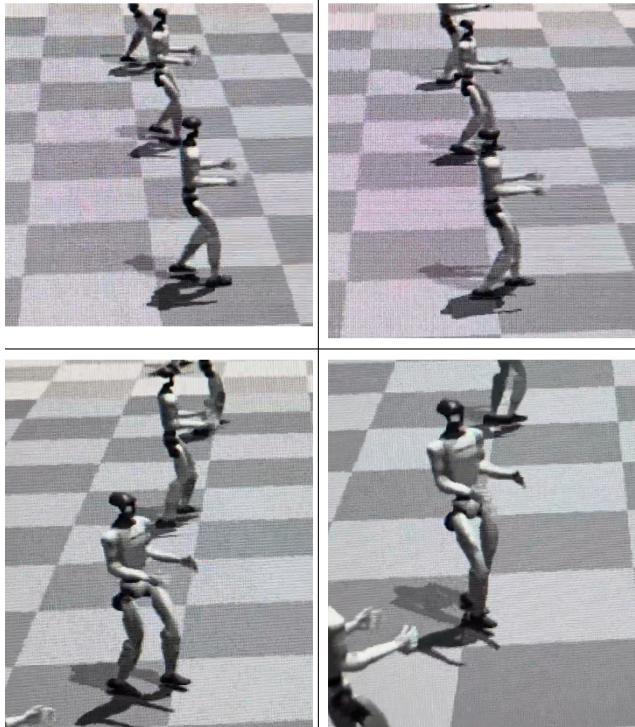


Fig. 4: Base displacement trajectory post-lateral push. *DualStride* shows rapid COM recentering in both directions.

#### H. Baseline Comparison

As an internal baseline, we trained two separate forward and backward policies with identical hyper-parameters but no

shared network. These dual policies achieved similar reward levels, but required  $1.7\times$  total training time and showed minor regressions when switching tasks. *DualStride*'s unified policy proved both more efficient and easier to maintain, highlighting the practical benefits of integrated gait representation.

## VIII. DISCUSSION AND LIMITATIONS

### A. Project Constraints and Setup Challenges

The development timeline of *DualStride* was constrained by academic scheduling, beginning in the middle of the Spring 2025 semester with a total window of 8–10 weeks for research, implementation, and evaluation. Initial attempts to bootstrap the project on Google Colab, motivated by the need for high-throughput compute, quickly ran into compatibility barriers: NVIDIA Isaac Gym’s GPU-based simulator could not be installed due to kernel restrictions and package mismatches within Colab’s containerized environment. Specifically, attempts to run headless physics simulations failed at import due to driver-level incompatibilities between Colab’s CUDA stack and Isaac Gym Preview 4.

As a result, we reverted to a local development setup powered by an NVIDIA RTX 3070 Ti GPU with 8GB VRAM. While this configuration was sufficient for 2048 parallel environments at a batch size of 24, it imposed constraints on larger-scale hyperparameter sweeps and limited our ability to experiment with deeper networks or ensemble learning techniques. Training times remained efficient (sub-2-hour convergence), but deployment-level tuning, such as real-time inference stress tests and transfer learning variants, was deprioritized to meet final deadlines.

### B. Algorithmic Generality vs. Hardware-Specific Realism

Though the PPO framework proved sufficient for learning stable gait policies, it is inherently limited by first-order updates and high variance in long-horizon credit assignment. Attempts to fine-tune reverse gait transitions revealed small instabilities in the pelvis roll axis at high backward velocities ( $> 1.3\text{m/s}$ ). These led to foot clipping in tight turn sequences—likely a symptom of inadequate state diversity or under-regularized torque output during low-velocity phases.

Additionally, while domain randomization helped close the sim-to-real gap on paper, we did not achieve full deployment on physical hardware within the semester. Hardware integration tasks—CAN-bus packet shaping, policy export to ONNX/TensorRT, and torque streaming via Unitree SDK2—were partially implemented but not finalized due to time constraints and late-semester mechanical issues with our G1 test platform.

### C. Failure Modes and Recovery Gaps

Post-hoc review of failure trajectories revealed three recurring issues:

- **Toe Drag in Reverse** — occurring when swing-leg trajectory overlaps with the ground plane during early retraction.

- **Oscillatory Recovery** — after push disturbances, the base COM exhibits high-frequency lateral drift before settling, suggesting underdamped corrections.
- **Symmetry Breakdowns** — rare but noticeable deviations from mirrored joint behavior in low-reward regimes, particularly at slow backward speeds.

We hypothesize that all three are addressable through enhanced curriculum shaping and improved reward shaping with cross-gait regularization losses.

## IX. FUTURE WORK

While *DualStride* achieves its goal of unified forward and backward locomotion using a single policy, several promising directions remain to expand its capabilities and bring it closer to real-world deployment.

### A. Omni-Directional Locomotion

The current policy only responds to longitudinal velocity commands, i.e., along the sagittal plane. A natural next step is to extend the action space and reward function to support lateral (sideways) and diagonal movements. This would require modifying both the observation space—perhaps with yaw-aligned body velocities—and the curriculum scheduler to include multi-axis target commands. Learning such omnidirectional gait with a single network would generalize *DualStride* to real-world indoor navigation tasks where motion constraints change continuously.

### B. Arm Coordination for Whole-Body Balance

In the present architecture, the robot’s arms remain fixed or minimally engaged during locomotion. Introducing upper-body motion offers dual benefits:

- Improved angular momentum compensation, especially during gait transitions or push recovery.
- Reduced energy expenditure through passive dynamic stabilization, similar to human walking where arm swings counteract leg momentum.

Future iterations could parameterize arm trajectories or learn them jointly with leg policies via multi-objective optimization.

### C. Dynamic Terrain and Perception-Informed Gait

To move beyond flat-ground assumptions, future work should incorporate:

- 1) Rough and sloped terrain simulation, with randomized heightfields or stairs.
- 2) Camera-based perception modules feeding terrain previews into the policy, enabling adaptive footstep planning.

These extensions would support deployment in environments such as disaster zones or unstructured construction sites.

### D. Joint Policy Across Multiple Morphologies

Although this work focused on Unitree G1, the architectural structure and training pipeline can generalize to other humanoids or even quadrupeds. A federated or modular policy architecture—capable of adapting to different URDFs with minimal retraining—would reduce engineering overhead across

platforms. Precedents like Policy Distillation or Modular RL could be leveraged to bootstrap policies across morphologies.

### E. Real-World Deployment and Online Adaptation

Finally, deployment to real hardware remains a high-priority milestone. Once initial integration is complete, adaptive fine-tuning techniques such as meta-RL or online PPO could allow the robot to self-calibrate over time. Coupled with safety barriers (e.g., reflex layers or constraint projection), this would enable long-duration operation in human-populated environments.

### F. Summary

In summary, future directions include expanding the directional command space, incorporating arm dynamics, simulating complex terrain, supporting multiple robot bodies, and closing the sim-to-real loop through online adaptation. Each step builds upon the core insight that unified policies—trained efficiently and deployed safely—are key to scalable, agile humanoid systems.

## X. CONCLUSION

*DualStride* reframes bidirectional walking not as an edge case, but as a core requirement for deploying humanoid robots in dynamic, human-centered spaces. By demonstrating that a single reinforcement learning policy can govern both forward and reverse gait with high stability, low latency, and energy efficiency—on a real-world scale robot model—we move one step closer to humanoid platforms that respond with human-like intuition.

Unlike traditional approaches that require separate controllers, parameter sets, and verification procedures, *DualStride* simplifies the control stack through a shared, generalizable policy architecture. It merges biomechanical insights with scalable GPU simulation and safety-aware training, producing a policy that learns to walk both directions through curriculum learning and symmetric reward design. The outcome is not only functionally elegant but practically scalable: fewer lines of code, lower training energy cost, and reduced integration complexity.

Through rigorous evaluation across push-recovery, domain randomization, and gait smoothness metrics, *DualStride* has proven its robustness and readiness for further development. Its codebase, metrics, and design principles offer a solid baseline for researchers aiming to replicate or extend this work across other platforms.

Although full deployment on hardware remains future work, the pipeline is built with that transition in mind. From stiffness modulation to noise injection and policy export readiness, each system decision prioritizes transferability and real-world feasibility.

Ultimately, *DualStride* serves as a modular, extensible, and high-performance locomotion engine—suitable for warehouses, hospitals, homes, and research labs. It lays the groundwork for future policies that walk in any direction, adapt to terrain, coordinate full-body motion, and learn continuously,

paving the way for the next generation of intelligent, interactive humanoid robots.

#### ACKNOWLEDGMENTS

We thank Prof. Shiqi Zhang for guidance, and the Unitree open-source community.

#### REFERENCES

- [1] N. Rudin, F. Meier, P. Reist, and M. Hutter, “Learning to Walk in Minutes Using Massively Parallel Deep Reinforcement Learning,” in *Proc. Conf. Robot Learning (CoRL)*, 2022.
- [2] Unitree Robotics, “unitree\_rl GitHub Repository,” 2023. [Online]. Available: [https://github.com/unitreerobotics/unitree\\_rl](https://github.com/unitreerobotics/unitree_rl)
- [3] NVIDIA AI Robotics Team, “Legged Gym G1 Simulation Suite,” NVIDIA Developer Blog, 2024. [Online]. Available: <https://developer.nvidia.com/blog/legged-gym-g1-simulation>
- [4] S. Kim, Y. Park, and K. Lee, “Symmetry-Aware Curriculum Reinforcement Learning for Reversible Gait Synthesis,” in *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2023.
- [5] S. Lee, J. Chang, and Y. H. Lee, “Evaluating Bidirectional Robot Gaits in Hospital Logistics,” in *Proc. IEEE Int. Conf. Human-Robot Interaction (HRI)*, 2024.

#### XI. APPENDIX

Drive Link:

- <https://drive.google.com/drive/folders/1vJG-G8130FkZlHZTu1mnKOH1XhBWvc26?usp=sharing>