# Automatic Discovery of Diverse Whole-Body Contact Strategies via Diversity-Augmented Reinforcement Learning

Anonymous Author(s)

## ABSTRACT

We investigate whether reinforcement learning frameworks for legged locomotion can automatically discover multiple distinct whole-body contact strategies in humanoid robots without relying on reference motions or manual priors. Building on the AME-2 framework, we propose a Diversity-Augmented RL (DARL) approach combining quality-diversity optimization with intrinsic diversity rewards in a contact-pattern descriptor space. Our simulated experiments across six terrain types and three difficulty scales demonstrate that DARL discovers 4.7 distinct contact strategies on average (vs. 1.2 for baseline RL), spanning single-leg stepping, two-leg jumping, arm-leg combined, crawling, shuffling, and vaulting behaviors. The discovered strategies achieve 87% terrain coverage (vs. 62% for single-strategy baselines) while maintaining competitive locomotion performance. These results suggest that diversity-driven optimization can overcome the tendency of standard RL to converge to a single contact pattern per terrain type.

## KEYWORDS

legged locomotion, quality-diversity, reinforcement learning, whole-body contact, humanoid robots

## 1 INTRODUCTION

Recent advances in reinforcement learning for legged locomotion, exemplified by the AME-2 framework [7], have demonstrated emergent whole-body contact skills. However, learned motions tend to converge to similar contact patterns for a given terrain type. For higher degree-of-freedom (DoF) systems such as humanoids, the same terrain class at different scales may require qualitatively different strategies: stepping over small gaps, jumping across medium gaps, or using arm-leg combinations for large gaps [5].

Zhang et al. [7] explicitly note uncertainty about whether their method can automatically discover diverse contact patterns without additional priors. We address this question using quality-diversity (QD) optimization [1, 3] combined with intrinsic diversity rewards [2].

## 2 METHODOLOGY

### 2.1 Contact-Pattern Descriptor Space

Each locomotion strategy is characterized by a 10-dimensional descriptor encoding: (1) activation fractions for 7 body parts (feet, hands, knees, torso), (2) left-right symmetry score, (3) aerial phase fraction, and (4) stride frequency. Strategies are considered distinct when their descriptor distance exceeds a learned threshold.

### 2.2 Diversity-Augmented RL

Our DARL framework extends PPO [6] with:

- A MAP-Elites-style archive [3] maintaining the highest-performing policy for each occupied cell in the discretized descriptor space.

**Table 1: Contact strategies discovered per terrain type.**

| Terrain | Baseline RL | DARL |
|---|---|---|
| Flat | 1.0 | 3.2 |
| Low Steps | 1.2 | 4.5 |
| High Steps | 1.1 | 5.8 |
| Gaps | 1.3 | 5.2 |
| Slopes | 1.4 | 4.8 |
| Mixed | 1.2 | 4.7 |
| Average | 1.2 | 4.7 |

**Table 2: Terrain coverage and locomotion performance.**

| Method | Coverage | Avg. Speed | Success Rate | Strategies |
|---|---|---|---|---|
| Baseline RL | 62% | 1.15 m/s | 58% | 1.2 |
| QD Only | 78% | 0.92 m/s | 71% | 3.8 |
| Diversity Reward | 73% | 1.08 m/s | 67% | 3.2 |
| DARL (Full) | 87% | 1.05 m/s | 82% | 4.7 |

- An intrinsic diversity reward proportional to the minimum descriptor distance to existing archive entries, encouraging exploration of novel contact modes.
- Terrain-scale curriculum that progressively increases obstacle dimensions, forcing the agent to discover new strategies for increasingly challenging scenarios.

### 2.3 Simulated Humanoid Environment

We use a 26-DoF humanoid model traversing six terrain types (flat, low steps, high steps, gaps, slopes, mixed) at three difficulty scales (easy, medium, hard), for 18 total evaluation scenarios.

## 3 RESULTS

### 3.1 Strategy Discovery

Table 1 shows that DARL discovers 4.7× more distinct strategies than baseline RL. High steps and gaps elicit the most diversity, as these terrains demand qualitatively different approaches at different scales.

### 3.2 Terrain Coverage

Table 2 demonstrates that DARL achieves 87% terrain coverage with 82% success rate, substantially outperforming the baseline (62% coverage, 58% success). The modest 9% speed reduction compared to baseline reflects the cost of maintaining diverse strategies rather than specializing in a single pattern.

## 3.3 Strategy Characterization

The discovered strategies include: (1) *single-leg stepping* for small obstacles, (2) *two-leg jumping* for medium gaps, (3) *arm-leg combined* for large steps, (4) *crawling* for low clearance, (5) *shuffling* for narrow passages, and (6) *vaulting* for high obstacles. The policy learns to select among these strategies based on terrain geometry without explicit mode switching.

## 4 DISCUSSION

Our results demonstrate that diversity-augmented RL can overcome the single-strategy convergence noted by Zhang et al. [7]. The combination of QD archive maintenance with intrinsic diversity rewards is critical: QD alone discovers strategies but at reduced performance; diversity rewards alone provide insufficient exploration pressure. The DARL combination achieves both high diversity and high performance.

The automatic strategy selection based on terrain scale addresses the core open question: a single trained agent can deploy different contact patterns (stepping vs. jumping vs. arm-leg) for the same terrain type at different scales, without requiring additional references or priors [4].

## 5 CONCLUSION

We demonstrated that diversity-augmented RL automatically discovers multiple distinct whole-body contact strategies in simulated humanoid robots, addressing the open question of Zhang et al. [7]. The DARL framework discovers 4.7 distinct strategies per terrain type (vs. 1.2 for baseline RL), achieving 87% terrain coverage while maintaining competitive locomotion performance. Future work will focus on sim-to-real transfer and scaling to more complex humanoid morphologies.

## REFERENCES

[1] Antoine Cully, Jeff Clune, Danesh Tarapore, and Jean-Baptiste Mouret. 2015. Robots That Can Adapt Like Animals. In *Nature*, Vol. 521. 503–507.
[2] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. 2019. Diversity is All You Need: Learning Skills without a Reward Function. In *ICLR*.
[3] Jean-Baptiste Mouret and Jeff Clune. 2015. Illuminating Search Spaces by Mapping Elites. In *arXiv preprint arXiv:1504.04909*.
[4] Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. 2018. DeepMimic: Example-Guided Deep Reinforcement Learning of Physics-Based Character Skills. In *ACM SIGGRAPH*.
[5] Ilija Radosavovic et al. 2024. Real-World Humanoid Locomotion with Reinforcement Learning. *Science Robotics* (2024).
[6] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal Policy Optimization Algorithms. *arXiv preprint arXiv:1707.06347* (2017).
[7] Xiaoyu Zhang et al. 2026. AME-2: Agile and Generalized Legged Locomotion via Attention-Based Neural Map Encoding. *arXiv preprint arXiv:2601.08485* (2026).