
Mimic under Adversarial Perturbations (MAP)

Anirudh Shrihari
Carnegie Mellon University
ashriha2@andrew.cmu.edu

Prajwal Gurunath
Carnegie Mellon University
pgurunat@andrew.cmu.edu

Manyung Emma Hon
Carnegie Mellon University
mehon@andrew.cmu.edu

Abstract

1 Humanoid motion imitation methods routinely demonstrate compelling behaviors
2 such as parkour, dancing, and running, yet these skills are typically trained for
3 exact reproduction and evaluated only in idealized settings without real-world
4 perturbations. Conventional reinforcement learning approaches rely on carefully
5 engineered rewards to induce natural gait patterns, making training highly experi-
6 mental and time-consuming. In this work, we study the robustness and practicality
7 of two adversarial imitation learning methods (1) Adversarial Motion Priors (AMP)
8 and (2) Adversarial Differential Discriminator (ADD), as a means to eliminate
9 reward tuning while enabling motion policies that remain reliable under external
10 perturbations. We further extend these methods to support steerability, a require-
11 ment for real-world deployment where human motion data must guide, rather
12 than dictate, robot behavior. We introduce a comprehensive evaluation framework
13 consisting of: (1) custom reference motions for the Unitree G1 humanoid that
14 simulate box-carrying postures, (2) curriculum-based force-perturbation training
15 using both constant and random external forces, and (3) an additional steering
16 objective for ADD. Our experiments show that ADD exhibits faster and more
17 stable learning, achieving significantly longer rollouts, while AMP produces higher
18 motion fidelity but suffers from reduced tracking accuracy under disturbance. We
19 discuss the trade-offs between these approaches and present strategies that enable
20 robust, steerable, and practically deployable humanoid motion policies.

21 1 Introduction

22 Physics-based character animation and motion imitation have made remarkable progress with deep
23 reinforcement learning [1, 2, 4]. Modern approaches can replicate complex human motions including
24 acrobatics, dance, and athletic movements with high fidelity. However, these methods (1) are
25 essentially replaying the motions and (2) are typically evaluated in pristine simulation environments
26 without the perturbations common in real-world scenarios.

27 Real-world robotic applications demand more than just motion reproduction; they require robustness
28 to external forces, payload handling, and environmental variations. A humanoid robot must maintain
29 stability and natural movement while carrying objects, responding to unexpected pushes, or operating
30 on varied terrain. This gap between impressive demonstrations and practical utility motivates our
31 investigation.

32 **Problem Formulation.** We focus on a specific but representative challenge: training humanoid
33 robots to accept directional commands and walk naturally while carrying a heavy item i.e., a box.
34 This task combines several difficulties: (1) continuous external forces on the end-effectors, (2) need
35 to maintain human-like gait patterns from reference motions, and (3) follow a steering command for
36 the locomotion. We would also like to do this with minimal reward tuning. To ensure we meet our
37 criteria we employ and evaluate two state-of-the-art adversarial imitation learning methods: AMP [2]
38 and ADD [4].

39 **Contributions.** Our work makes the following contributions:

- 40 • Motion retargeting from popular human motion datasets such as AMASS [9] and LAFAN
41 [10] to the Unitree G1
- 42 • Custom dataset generation for box-carrying reference motions on the Unitree G1 humanoid
- 43 • A systematic evaluation of AMP and ADD under random external force perturbations,
44 including large constant forces (for box carrying) and smaller random impulses (for enhanced
45 robustness)
- 46 • Curriculum learning framework to ensure model training stability
- 47 • Implementation of steering rewards (base velocity and heading) for joint optimization of
48 task objectives and motion imitation
- 49 • Empirical analysis revealing distinct failure modes and trade-offs: ADD’s speed vs. natural-
50 ness, AMP’s quality vs. tracking accuracy

51 2 Related Work

52 **Physics-Based Character Control.** DeepMimic [1] pioneered using deep RL for physics-based
53 character animation by combining imitation rewards with task rewards. AMP [2] introduced adversar-
54 ial motion priors, using a discriminator to distinguish policy motion from reference motion without
55 explicit pose tracking rewards. The emphasis was on reproducing the "style" of the motions.

56 **Differential Discriminators.** ADD [4] proposed using differential observations (relative differences
57 between policy and reference) rather than absolute states for the discriminator. This approach enables
58 the policy to track motion patterns while adapting to external forces or different morphologies.

59 **Robustness in Motion Imitation.** While several works address domain randomization for sim-to-real
60 transfer [8], systematic studies of external force robustness in adversarial imitation learning remain
61 limited. FALCON [6] employs a feasible force curriculum to train humanoids for force-adaptive
62 manipulation, but relies on extensive reward engineering rather than leveraging motion priors for
63 natural gait patterns.

64 Our work addresses this gap by systematically evaluating adversarial motion imitation methods (AMP
65 and ADD) under external force perturbations while incorporating steering commands. This approach
66 combines natural motion quality from reference data with robustness to real-world forces, reducing
67 the need for manual reward tuning.

68 3 Methods

69 3.1 Background: AMP and ADD

70 AMP uses adversarial imitation learning where a discriminator D tries to distinguish between policy-
71 generated motions and reference motions, while the policy learns to fool the discriminator. The
72 discriminator is trained via adversarial loss while the policy is trained via reinforcement learning with
73 discriminator-based rewards. ADD also uses a discriminator D , but the input to the discriminator is
74 the delta of various parameters of the policy and the reference motion. The discriminator provides a
75 high reward if the policy motion is close to or the same as the reference motion.

76 **AMP** [2] uses absolute state observations for discrimination. The discriminator $D_\theta : \mathcal{S} \rightarrow [0, 1]$ is
77 trained to distinguish policy states from reference states using binary cross-entropy loss. The style
78 reward is computed as:

$$r_{style}(s_t) = -\log(1 - D_\theta(s_t)) \quad (1)$$

79 where s_t contains joint positions, velocities, and root information from the policy at time t . This
80 reward encourages the policy to generate states that fool the discriminator (achieving $D_\theta(s_t) \rightarrow 1$).

81 **ADD** [4] uses differential observations computed as the difference between policy and reference
82 states. The discriminator observes:

$$\mathbf{o}_t^{diff} = \phi(s_t^\pi, s_t^{ref}) \quad (2)$$

83 where ϕ computes relative features including root position difference $\Delta \mathbf{p}_{root}$, root velocity difference
84 $\Delta \mathbf{v}_{root}$, and relative body positions with respect to the root. The discriminator $D_\theta : \mathcal{O}^{diff} \rightarrow [0, 1]$

85 is trained with the zero vector $\mathbf{0}$ as the only positive sample (representing perfect tracking) and
 86 error vectors \mathbf{o}_t^{diff} from the policy as negative samples. The reward encourages minimizing the
 87 differential:

$$r_{style}(s_t^\pi, s_t^{ref}) = -\log(1 - D_\theta(\mathbf{o}_t^{diff})) \quad (3)$$

88 This differential formulation provides translation and rotation invariance, allowing policies to adapt
 89 their global position while maintaining relative motion patterns.

90 3.2 Force Perturbation Framework

91 We extended the MimicKit codebase to support configurable force perturbations:

92 **Random Force Model.** Forces are applied probabilistically:

$$\mathbf{f}_t = \mathbf{f}_{t-1} + \mathbf{f}_{new} \cdot \mathbb{I}(u < p_{force}) \quad (4)$$

$$\mathbf{f}_{new} \sim \mathcal{U}(-1, 1)^3 \odot \mathbf{s}_{force} \quad (5)$$

93 where p_{force} is application probability (0.2 per step), \mathbf{s}_{force} is the force scale vector, and $u \sim \mathcal{U}(0, 1)$.
 94 The indicator function $\mathbb{I}(\cdot)$ ensures forces are applied stochastically.

95 **Curriculum Learning.** We gradually increase force magnitudes during training using a linear ramp:

$$\mathbf{s}_{force}(t) = \begin{cases} \mathbf{0} & \text{if } t < 3000 \\ \min(\mathbf{s}_{max}, \frac{t-3000}{6000} \cdot \mathbf{s}_{max}) & \text{otherwise} \end{cases} \quad (6)$$

96 where t is the training iteration and \mathbf{s}_{max} is the target maximum force vector. Forces remain at zero
 97 for the first 3000 iterations to allow initial policy learning, then linearly increase to \mathbf{s}_{max} over the next
 98 6000 iterations, reaching maximum force at iteration 9000. The min operator is applied element-wise
 99 to cap each force component at its maximum value.

100 3.3 Box-Carrying Dataset Generation

101 We created reference motions for box-carrying by:

- 102 1. Retargeted human walking motion from AMASS [9] and LAFAN [10] to Unitree G1
- 103 2. Manually filtered out walking gaits with forward and changing direction
- 104 3. Initialize shoulder, elbow, and wrist joints to a fixed box-holding pose (elbow to wrist
 105 extended forward)

106 This generates naturalistic reference data where the lower body exhibits natural walking motion while
 107 the upper body maintains a carrying pose. The custom initialization of the upper-body joints ensure
 108 the reference motion represents realistic box-carrying behavior.

109 3.4 Steering Task with ADD

110 Given AMP’s position drift issues under force (Section 5) and slower training to ensure stability, we
 111 developed a steerable policy using ADD to enable directional control while maintaining motion style.
 112 We implemented a steering environment where the robot must follow dynamically changing target
 113 directions while maintaining natural locomotion patterns.

114 **Dataset Preparation.** We retargeted the LAFAN dataset [10] to the Unitree G1, using 6 diverse
 115 motion sets with neutral (natural) upper-body positions. This provides varied locomotion patterns
 116 including walking, running, and transitions.

117 **Discriminator Enhancement.** We augmented ADD’s differential discriminator observations with
 118 velocity information, enabling the discriminator to assess motion quality in the context of current
 119 movement direction.

120 **Task Reward.** The steering task reward combines directional objectives:

$$r_{task} = w_{tar} \cdot r_{vel} + w_{face} \cdot r_{face} \quad (7)$$

121 where $r_{vel} = \exp(-\|v - v_{tar}\mathbf{d}_{tar}\|^2)$ rewards moving at target speed v_{tar} in target direction \mathbf{d}_{tar} ,
 122 and $r_{face} = \exp(-\|\mathbf{d}_{heading} - \mathbf{d}_{tar}\|^2)$ rewards facing the target direction. The total reward
 123 combines task and style objectives:

$$r_{total} = r_{task} + w_{style} \cdot r_{style} \quad (8)$$

124 The choice of ADD over AMP for steering was motivated by ADD’s superior position tracking under
 125 force, making it more suitable for directional control tasks despite overcompensation tendencies in
 126 single-axis force training.

127 4 Experimental Setup

128 **Environment.** We use Isaac Gym [7] with the Unitree G1 humanoid (29 DoF) on flat terrain. Control
 129 frequency is 30Hz with simulation at 120Hz. Episode length is 10 seconds with early termination on
 130 falls.

131 **Training Details.** We use PPO [5] with 4096 parallel environments. Both AMP and ADD use
 132 identical network architectures and base hyperparameters, with key differences only in discriminator
 133 formulation and batch size. Networks are 2-layer fully-connected with ReLU activation: 1024 units
 134 in the first layer and 512 in the second. Training uses SGD optimizer with learning rate 1×10^{-4}
 135 and gradient clipping (max norm 1.0). PPO hyperparameters include: discount factor $\gamma=0.99$, GAE
 136 $\lambda=0.95$, clip ratio 0.2, 5 update epochs per iteration, and 32 steps per iteration. For AMP, we use batch
 137 size 4 and reduce discriminator strength (loss weight 2.0, gradient penalty 0.5) to prevent discriminator
 138 dominance. ADD uses batch size 12. Both methods use discriminator logit regularization 0.01 and
 139 weight decay 0.0001. Training runs for 20,000-25,000 iterations (2-3 billion samples, 15-20 GPU
 140 hours on an RTX 4090).

141 **Force Curriculum.** Forces gradually increase from 0 to maximum magnitude over 192,000 envi-
 142 ronment steps (6000 iterations) starting after 3000 global steps. For [10N, 10N, 30N] configuration,
 143 forces start at [0,0,0] and linearly increase to [10N, 10N, 30N]. Forces are applied to left and right
 144 wrist links, remain constant within each episode, but maximum magnitude increases across episodes
 145 as the curriculum progresses.

146 **Force Configurations.** We test these force scenarios:

- 147 • **Constant Light:** 30N downward force on wrists (simulating 3kg box)
- 148 • **Constant Heavy:** 100N downward force on wrists (simulating 10kg box)
- 149 • **Multi-Directional:** [10N, 10N, 30N] forces in x, y, z axes to reduce single-axis overcom-
 150 pensation

151 **Evaluation Metrics.** We measure:

- 152 • Episode length (surrogate for stability and training progress)
- 153 • Pose tracking error: Mean joint-wise position error from reference motion
- 154 • Root tracking error: Position and velocity deviation from reference trajectory
- 155 • End-effector position error (feet and hands)
- 156 • Discriminator reward (style matching quality)

157 5 Results

158 We present our results: (1) force perturbation training on end effectors vs baselines box carrying gait
 159 with forces at test time, (2) qualitative analysis revealing the overcompensation problem, and (3)
 160 quantitative and qualitative analysis of steerable policy with ADD.

161 5.1 Force Perturbation Training:

162 We first establish a baseline by training both the ADD and AMP methods using a single walking gait
 163 reference motion, without any force-based curriculum. We then introduce force perturbation training

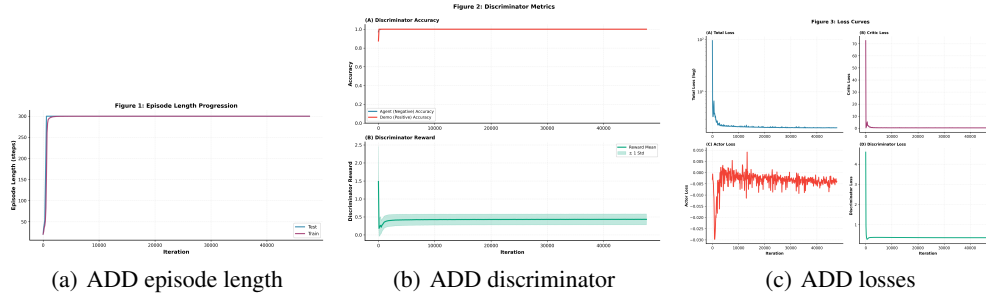


Figure 1: **ADD training curves under 30N constant force.** Episode length saturates quickly (a), discriminator loss remains high (b), but policy converges (c).

164 using a curriculum of external forces with magnitudes $[10\text{ N}, 10\text{ N}, 30\text{ N}]$, applied to the same walking
 165 gait. During training, we apply external forces and compare performance across multiple metrics,
 166 including velocity and position tracking, to assess the ability of the policies to counteract disturbances
 167 while preserving smooth and natural gait dynamics. Importantly, no additional force-specific reward
 168 terms are introduced during training. We also qualitatively compare the upper-body position tracking
 169 and gait at test time.

170 **ADD Performance.** Figure 1 shows ADD training curves. Key observations: (1) Episode length
 171 saturates quickly (within 5000 iterations), reaching maximum rollout duration, (2) Discriminator
 172 maintains high loss, suggesting policy motion differs from reference despite achieving stability (3)
 173 Policy learns rapid force compensation, enabling continued walking

174 **AMP Performance.** Figure 2 shows AMP training after hyperparameter tuning (gradient clipping,
 175 discriminator weakening): (1) More gradual learning curve than ADD, reaching stability around
 176 15,000 iterations, (2) Discriminator reward steadily increases, indicating better style matching (3)
 177 Higher pose tracking accuracy but lower stability under force

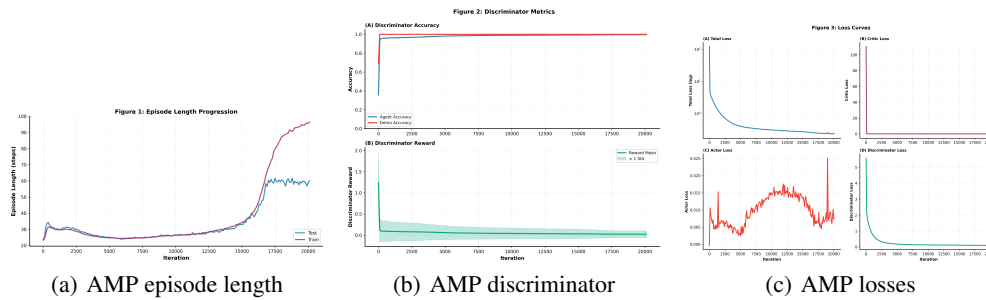


Figure 2: **AMP training curves under force perturbations.** More gradual learning (a), improving discriminator reward (b), and stable loss convergence (c).

178 5.2 Qualitative Analysis: The Overcompensation Problem

179 Having established baseline training, we now investigate policy behavior under different force
 180 magnitudes to understand robustness limits.

181 **Single-Axis Forces Cause Overcompensation in ADD** Under $[0,0,30\text{N}]$ curriculum training, ADD
 182 produces stable walking but develops unnatural compensation strategies: (1) Arms raised above
 183 head to counterbalance downward forces (2) Exaggerated leaning motions during gait (3) Persistent
 184 overcompensation: Even at zero force, policy maintains defensive posture.

185 The $[0,0,100\text{N}]$ curriculum exacerbates this problem. Figure 3 shows that the aggressive curriculum
 186 creates policies that overcompensate severely even without applied forces, sometimes becoming

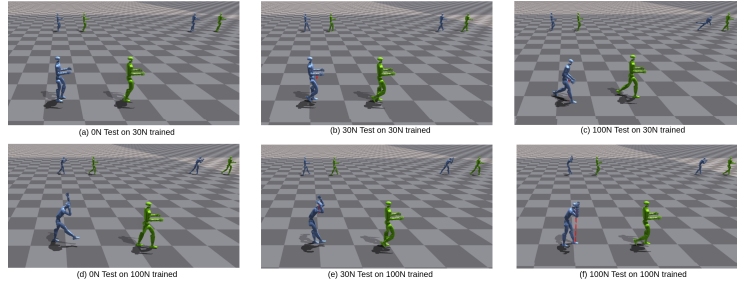


Figure 3: **ADD behavior under single-axis force curricula.** 30N curriculum allows adaptation but shows arm raising, while 100N curriculum causes severe overcompensation even at zero force.

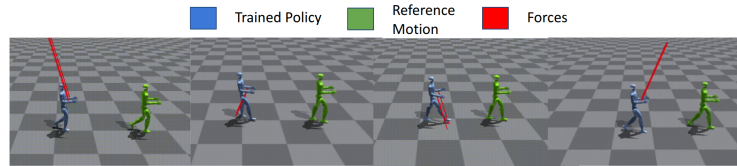


Figure 4: **Effect of multi-directional forces on ADD.** Multi-directional [10N,10N,30N] training produces natural poses.

187 unstable. This behavior suggests the policy learns axis-specific compensation strategies (raising arms
 188 for downward forces) that persist inappropriately.

189 **Multi-Directional Forces Solve Overcompensation** To address this issue, we trained ADD with
 190 [10N, 10N, 30N] multi-directional forces. This configuration applies forces in x, y, and z axes,
 191 preventing axis-specific compensation strategies. Figure 4 shows that multi-directional training
 192 produces natural poses even after force removal, validating our hypothesis that force diversity
 193 prevents overcompensation.

194 **AMP: Better Pose Tracking, Lower Force Robustness** AMP maintains more natural-looking
 195 poses and better pose tracking accuracy but exhibits lower force robustness. We trained AMP with
 196 the same three force configurations as ADD to enable direct comparison: (1) Significant drift in
 197 root position during walking under force, (2) Occasional stumbling and falls under moderate forces
 198 (>40N) (3) Better recovery when forces are removed compared to ADD, (4) Maintains closer match
 199 to reference motion (higher discriminator rewards)

200 Multi-directional forces: (1) Improved stability compared to single-axis training, (2) Still exhibits
 201 position drift, making it less suitable for directional control, (3) Maintains natural pose quality across
 202 all force magnitudes. Figure 5 shows that AMP can walk with light constant forces but experiences
 203 drift and instability with strong perturbations, confirming the trade-off between pose tracking accuracy
 204 and force robustness across all training configurations.

205 5.3 Quantitative Comparison

206 We now quantify the performance differences across force configurations and methods.

207 **Pose Tracking and Stability Metrics** Table 1 presents pose tracking error and episode length
 208 across force magnitudes for both AMP and ADD trained with all three force configurations. Key
 209 findings:

- 210 • **AMP:** Excellent pose tracking at 0N across all configs, but rapid degradation under force;
 211 multi-directional training provides modest improvement
- 212 • **ADD:** Lower pose tracking error under force but higher error at 0N with single-axis training
 213 (overcompensation)

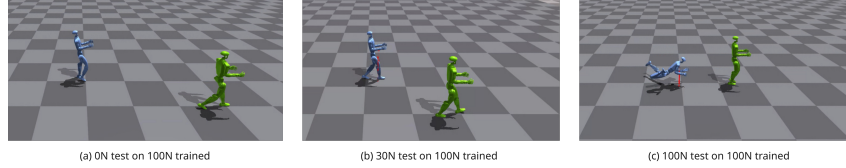


Figure 5: **AMP behavior under constant forces.** Model trained with curriculum up to 100N shows reasonable walking with light forces but drifts and falls under strong forces, suggesting a lower upper-bound when compared to ADD.

- 214 • **Multi-directional benefit:** [10,10,30] training reduces ADD’s 0N pose error while main-
215 taining robustness; AMP shows smaller improvement but still drifts.
- 216 • **Body Position Error:** ADD maintains a lower body position error with multi-directional
217 forces due to explicitly learning and weighting the difference between the reference motion,
218 hence the signal for tracking is more nuanced than AMP which instead learns the style of
219 the motions
- 220 • **Episode Length:** Given that the episode ends right when the robot falls, a higher episode
221 length results in continuous tracking of repeatable motions while also countering varying
222 forces during the episode. Multi-directional forces does allow the agent to be aware of
223 any incoming force to compensate and continue tracking positional errors and style of the
224 motions

Table 1: Training results under different force curricula (final metrics after convergence)

Force Curriculum	Episode Length (s)	Body Pos Error	Iterations
ADD [0,0,30]	19.2	0.211	15000
ADD [0,0,100]	15.6	0.291	15000
ADD [10,10,30]	260.5	0.015	15000
AMP [0,0,30]	31.2	0.107	15000
AMP [0,0,100]	24.9	0.085	15000
AMP [10,10,30]	265.2	0.045	15000

225 **Method Comparison** Table 1 summarizes the key differences between AMP and ADD under force
226 perturbations. ADD was also evaluated on steering tasks, where it demonstrated the ability to follow
227 target directions while maintaining learned motion styles.

228 5.4 Qualitative Analysis: Steerable Policy

229 We now measure the Qualitative performance of the Steerable Policy. Figure 6 show that the policy
230 can track the velocity, but it struggles with tracking the reference motion gait.

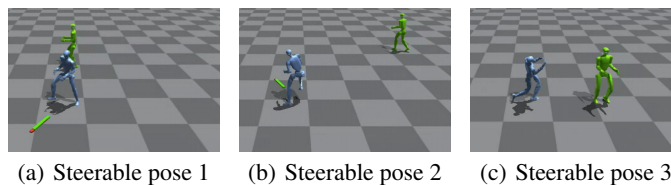


Figure 6: **Robot pose from steerable policy compared to reference motion**

231 From Figure 7, we can see that the body and root rotation error is pretty high, indicating that the
232 policy isn’t able to track the reference policy pretty well. The root velocity error is 0.4 m/s indicating
233 that it can track velocity pretty well.

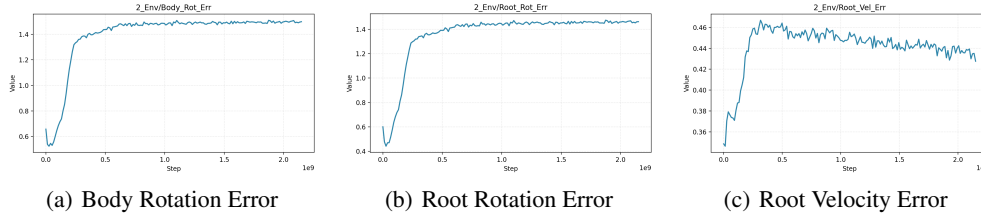


Figure 7: ADD Steerable training curves

234 6 Discussion

235 6.1 Why AMP Struggles with Pose Tracking Under Force

236 AMP’s absolute discriminator observes full state s , requiring close pose tracking to reference for high
 237 rewards. Under external forces:

- 238 1. The policy must balance discriminator reward (match reference pose) with task reward
 239 (avoid falling under force)
- 240 2. Forces cause accumulating root position errors that degrade overall pose tracking accuracy
- 241 3. Gradient clipping and discriminator weakening (necessary for training stability) reduce
 242 corrective signals for pose tracking

243 The trade-off is excellent pose tracking at 0N but rapid degradation under force. AMP maintains
 244 natural motion quality but cannot sustain pose tracking accuracy under sustained perturbations.

245 6.2 Curriculum Learning Insights

246 Our force curriculum reveals important insights about force magnitude and directionality:

247 Force Magnitude Effects:

- 248 • **30N curriculum:** Produces policies that handle 0-40N robustly without overcompensation
- 249 • **100N curriculum:** 100N is well past the carrying capacity of the G1 and hence it tries and
 250 overcompensates it’s motion

251 6.3 Why ADD isn’t able to walk with a smooth gait while tracking velocities

252 The output from ADD’s discriminator is very small while training. The velocity tracking reward
 253 outweighs the discriminator rewards, hence the ADD policy isn’t able to learn reference motion while
 254 also tryin to track velocity. This can be fixed by weighing the rewards from ADD higher than velocity
 255 tracking, or adding a curriculum to introduce velocity tracking after the policy has learnt how to
 256 imitate the reference walking motion.

257 6.4 Limitations and Future Work

258 Current Limitations:

- 259 • Evaluation limited to downward wrist forces; real box-carrying involves distributed forces
- 260 • Flat terrain only; uneven ground would compound challenges
- 261 • No real hardware validation; sim-to-real gap remains unquantified
- 262 • Limited analysis of failure modes (what specific forces cause falls?)

263 Future Directions:

- 264 • **Explicit pose constraints:** Add constraints on unnatural joint configurations

- 265 • **Multi-task training:** Jointly train on steering, force perturbations, and terrain variations
- 266 • **Real robot deployment:** Evaluate on physical Unitree G1 with actual box-carrying tasks
- 267 • **Improved steerable policy:** Improve the baseline steerable ADD policy to better follow
- 268 reference policy motion
- 269 • **Steerable with force compensation:** Add force curriculum to the steerable ADD policy to
- 270 walk and change directions with external force perturbations.

271 7 Conclusion

272 We studied whether adversarial motion imitation can produce humanoid locomotion policies that are
273 not only visually plausible, but also practical under real-world disturbances. Using the Unitree G1 in
274 IsaacGym, we compared AMP and ADD under curriculum-based external force perturbations applied
275 at the wrists to emulate box-carrying loads, and we further extended ADD with steering objectives to
276 support velocity tracking control.

277 Across the force perturbation experiments, ADD consistently learned faster and achieved longer stable
278 rollouts, indicating stronger disturbance tolerance and easier optimization. However, single-axis
279 force curricula caused: ADD can learn brittle, axis-specific compensation strategies even when forces
280 are removed. We showed that introducing multi-directional forces mitigates this overcompensation,
281 improving zero-force posture while preserving robustness. In contrast, AMP tended to produce
282 higher-fidelity motion and better pose matching in nominal conditions, but suffered from root drift
283 and reduced stability as force magnitudes increased, highlighting a fundamental trade-off between
284 strict style adherence and force robustness when the discriminator observes absolute states.

285 For steerability, our ADD-based policy successfully tracked commanded velocity targets, but exhibited
286 large body and root rotation errors relative to the reference motion, suggesting that the task reward
287 can dominate the style objective. Qualitatively there is a clear tracking of directional commands, but
288 the style is diluted. This points to a clear path forward: improve reward balancing (or schedule it
289 via curriculum), strengthen the style signal during early training, and only later introduce aggressive
290 steering demands.

291 Overall, our results suggest that adversarial motion priors can substantially reduce reward engineering
292 while enabling robust humanoid behaviors, but practical deployment requires careful perturbation
293 design and objective scheduling. Future work will extend perturbations beyond wrist-only loads,
294 incorporate uneven terrain and richer disturbance distributions, and validate sim-to-real transfer on
295 hardware. We also plan to unify steering and force-robust training into a single multi-task framework
296 to produce humanoid policies that are simultaneously natural, commandable, and resilient in realistic
297 operating conditions.

298 Acknowledgments

299 We thank the MimicKit authors for releasing their codebase and the course staff for valuable feedback.
300 Computational resources were provided by Carnegie Mellon University.

301 **References**

- 302 [1] Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. Deepmimic: Example-
303 guided deep reinforcement learning of physics-based character skills. *ACM Transactions on*
304 *Graphics*, 37(4):1–14, 2018.
- 305 [2] Xue Bin Peng, Ze Ma, Pieter Abbeel, Sergey Levine, and Angjoo Kanazawa. Amp: Adversarial
306 motion priors for stylized physics-based character control. *ACM Transactions on Graphics*,
307 40(4):1–15, 2021.
- 308 [3] Xue Bin Peng, Yunrong Guo, Lina Halper, Sergey Levine, and Sanja Fidler. Ase: Large-scale
309 reusable adversarial skill embeddings for physically simulated characters. *ACM Transactions*
310 *on Graphics*, 41(4):1–16, 2022.
- 311 [4] Ziyu Zhang, Sergey Bashkirov, Dun Yang, Yi Shi, Michael Taylor, and Xue Bin Peng. Physics-
312 based motion imitation with adversarial differential discriminators. In *SIGGRAPH Asia 2025*
313 *Conference Papers*, 2025.
- 314 [5] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
315 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 316 [6] Yuanhang Zhang, Yifu Yuan, Prajwal Gurunath, Ishita Gupta, Shayegan Omidshafiei, Ali-
317 akbar Agha-mohammadi, Marcell Vazquez-Chanlatte, Liam Pedersen, Tairan He, and Guanya
318 Shi. FALCON: Learning force-adaptive humanoid loco-manipulation. *arXiv preprint*
319 *arXiv:2505.06776*, 2025.
- 320 [7] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles
321 Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, and others. Isaac
322 gym: High performance GPU-based physics simulation for robot learning. *arXiv preprint*
323 *arXiv:2108.10470*, 2021.
- 324 [8] Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real
325 transfer of robotic control with dynamics randomization. In *IEEE International Conference on*
326 *Robotics and Automation*, pages 3803–3810, 2018.
- 327 [9] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black.
328 AMASS: Archive of motion capture as surface shapes. In *IEEE International Conference on*
329 *Computer Vision (ICCV)*, pages 5442–5451, October 2019.
- 330 [10] Félix G. Harvey, Mike Yurick, Derek Nowrouzezahrai, and Christopher Pal. Robust motion
331 in-betweening. *ACM Transactions on Graphics*, 39(4), 2020.