

Push-and-Step: From RL-Based Balance Recovery to Physical Simulation of Dense Crowds

Anonymous CVPR submission

Paper ID



Figure 1. Evolution of a dense crowd as a push propagates through physically simulated agents maintaining balance via steps and contacts.

001 1. Introduction

002 Crowd simulation has traditionally modeled collective be-
003 havior through local interaction rules in moderately dense
004 environments, representing individuals as simplified 2D ge-
005 ometries such as disks or particles [5, 8]. This paradigm
006 suffices when interactions remain social and non-contact.
007 In highly dense settings—packed subways, concert pits,
008 or emergency evacuations—interactions become predomi-
009 nantly physical: contact is unavoidable, and a character re-
010 covering balance inevitably exerts force on neighbors, prop-
011 agating and potentially amplifying disturbances through the
012 crowd. Falls, trampling, and crush injuries are real-world
013 consequences that 2D representations are ill-equipped to
014 model [16, 20].

015 We investigate full-body, physics-based simulation of
016 humanoid agents maintaining balance under external per-
017 turbations in dense crowds. The core challenge is three-
018 fold: (i) physically simulated bipeds are notoriously diffi-
019 cult to control under large impulsive loads; (ii) in multi-
020 agent settings, every attempted recovery can trigger sec-
021 ondary pushes; and (iii) social norms constrain how agents
022 interact—brief shoulder contacts are acceptable, full-body
023 collisions are not [6, 15].

024 To address this, we employ reinforcement learning (RL)
025 to train a balance-recovery policy that allows simulated
026 characters to maintain or regain balance through socially
027 aware interactions when subjected to unpredictable external
028 pushes. Our training pipeline consists of two stages: a sin-
029 gle character first learns isolated balance recovery via mo-
030 tion imitation, before being adapted to a multi-agent, dense-
031 crowd setting.

The code is available on [GitHub](#), and the main contribu-
tions are:

- 032 • **RL for omnidirectional balance recovery.** We intro- 033
034 duce an RL framework combining adversarial imitation 035
036 learning with physics-based balance rewards, enabling 037
038 agents to generalize from a compact reference dataset to 039
040 pushes of varying strength, direction, and duration. 041
042
- 043 • **Socially aware multi-agent interactions.** An online 044
045 hand-contact heuristic guides agents toward neighbors’ 046
047 shoulders for mechanically efficient and socially plausi- 048
049 ble balance dissipation, enabling a single-character policy 050
051 to extend naturally to multi-agent crowds. 052
- 053 • **Scalable dense-crowd simulation.** The learned policy 054
055 scales to large numbers of physically simulated agents, 056
057 opening new avenues for studying safety and collective 058
059 behavior in high-density settings. 060

061 2. Related Work

062 **Dense crowd simulation.** Prior work has focused on nav- 063
064 igation and social behaviors, using particle-based 2D mod- 065
066 els that can capture push-propagation waves but miss the 067
068 physical dynamics of limb-level interactions [8, 11, 13, 20]. 069
070 To our knowledge, full-body physical simulation of dense 071
072 crowds has not been previously attempted. 073

074 **Human balance strategies.** Maintaining bipedal balance 075
076 relies on ankle and hip torques for small perturbations, with 077
078 stepping to expand the Base of Support (BoS) for larger 079
080 ones [2, 3]. In dense environments, hands can further as- 081
082 sist recovery by exerting forces on neighbors [3, 6]. Social 083
084 context shapes these interactions: pushing is tolerated when 085
086 necessary for stability but must be directed at socially ap- 087
088 propriate body locations such as the shoulder [1, 15, 19]. 089

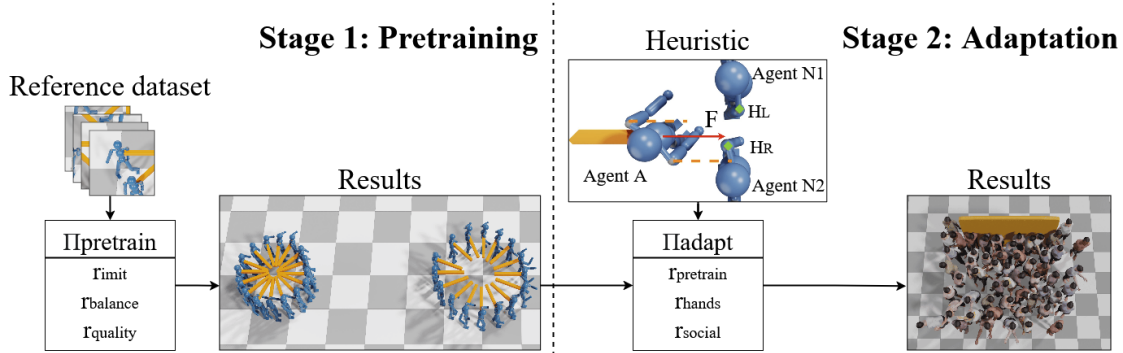


Figure 2. Two-stage training: pretraining (π_{pretrain}) for omnidirectional stepping, and adaptation (π_{adapt}) for multi-agent hand-shoulder interaction guided by the online heuristic.

063 **Physics-based character control.** Expert-designed
064 heuristics [4, 17, 26] and trajectory optimization [22, 27]
065 have been the traditional tools for physics-based balance.
066 Learning-based methods have since demonstrated strong
067 motion imitation [12, 23] and adaptation [10, 25]. Multi-
068 agent physics-based scenarios have been explored for
069 sports and locomotion [9, 21, 24], but dense physical crowd
070 interaction with balance recovery remains unaddressed.

071 3. Method

072 3.1. Overview

073 All agents are physically simulated humanoids controlled
074 via proportional-derivative (PD) servos [18]. The control
075 policy $\pi(\mathbf{a}_t | \mathbf{s}_t)$ takes the agent’s state \mathbf{s}_t (joint poses and
076 velocities over the past four frames) as input and outputs
077 target joint angles \mathbf{a}_t . Training proceeds in two stages (Fig-
078 ure 2): pretraining (π_{pretrain}) for isolated balance recovery,
079 and adaptation (π_{adapt}) for multi-agent interactions.

080 3.2. Stage 1: Adaptive Stepping Response

081 The first stage produces a policy that mimics human step-
082 ping strategies in response to external perturbations.

083 **Balancing dataset.** We collected a motion capture dataset
084 of a single subject pushed horizontally on the upper back
085 at eight evenly spaced angles. A wall placed 1 m ahead
086 prompted natural arm-raising motions, which are important
087 for the adaptation stage. Pushes range from 70 N to 200 N,
088 matching experimental observations [7].

089 **Training reward.** We follow the GAN-like ICCGAN
090 framework [23] with PPO [14] and an overall reward:

$$091 r_{\text{pretrain}} = w_i \frac{1}{2} (r_{\text{imit}} + 1) + w_g r_{\text{balance}} + w_q r_{\text{quality}}, \quad (1)$$

092 with weights $w_i=0.6$, $w_g=0.2$, $w_q=0.2$. The imitation re-
093 ward r_{imit} is the discriminator output (hinge loss in $[-1, 1]$)
094 measuring similarity to reference motions. The balance re-
095 ward penalizes deviations from physically stable postures:
096

$$097 r_{\text{balance}} = e^{-\|\text{CoM} - \text{CoM}_{\text{target}}\|} + e^{-\|\text{CoP} - \text{CoP}_{\text{target}}\|}. \quad (2)$$

098 $\text{CoM}_{\text{target}}$ is the center of the Base of Support at resting
099 height; $\text{CoP}_{\text{target}}$ is derived from momentum regula-
100 tion [22], shifting from ankle pressure adjustment to foot
101 placement when a step is needed:

$$102 \text{CoP}_{\text{target},x} = \text{CoM}_x + \frac{d_l p_x}{f_z} \text{CoM}_z + \frac{d_h L_y}{f_z}, \quad (3)$$

103 with damping factors $d_l=4$, $d_h=6$ and an analogous ex-
104 pression for the y component. The quality reward r_{quality}
105 suppresses artifacts when generalizing beyond the reference
106 dataset: it penalizes foot sliding (r_{foot}), heading deviation
107 (r_{heading}), and unnecessary kinetic energy (r_{effort}). During
108 training, push angles, magnitudes (70–200 N), and dura-
109 tions (0.7–1.3 s) are sampled uniformly, and small noise is
110 added to initial joint configurations to improve robustness.

111 3.3. Stage 2: Adaptation to Multi-Agent Interaction

112 The second stage adapts π_{pretrain} to multi-agent settings us-
113 ing the AdaptNet architecture [25], which freezes the pre-
114 trained policy and learns new contact behaviors through
115 latent-space injection and parallel adapter layers.

116 **Hand-contact heuristic.** Each agent’s hand targets are
117 computed online by an algorithm that considers candidate
118 shoulders of neighbors within 5 m. Targets are appended to
119 the state vector \mathbf{s}_t , and two objectives drive the selection:

- 120 1. *Collision avoidance:* the agent’s shoulder trajectory over
121 a one-second horizon (estimated from current linear mo-
122 mentum) is compared to all candidate shoulders. Any
123 shoulder within $\delta=0.25$ m of this trajectory is flagged as
124 a collision risk and selected as a hand target.
- 125 2. *Fall prevention:* for any hand without a collision-risk tar-
126 get, reachable shoulders (distance to the predicted po-
127 sition in $[0.1, 0.6]$ m) are evaluated by collinearity with
128 the agent’s linear momentum [19], and the most energy-
129 efficient one is selected.

130 **Adaptation reward.** The final reward extends the pre-
131 training reward with two new terms:

$$132 r_{\text{adapt}} = w_p r_{\text{pretrain}} + w_h r_{\text{hands}} + w_s r_{\text{social}}, \quad (4)$$

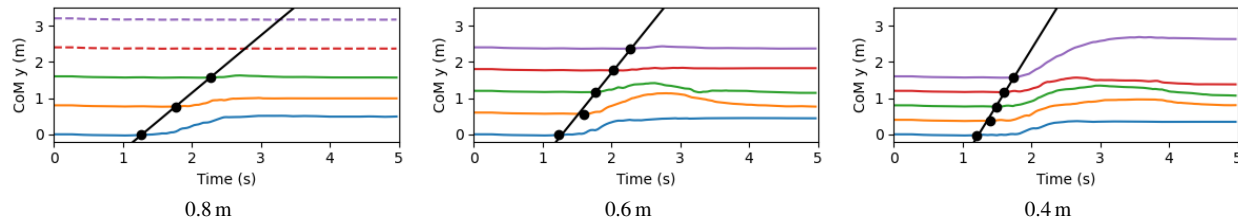


Figure 3. CoM trajectories in the line experiment for three interpersonal distances. Dashed lines indicate agents not yet reached by the push; dots mark contact events.

133 with $w_p=0.5$, $w_h=0.2$, $w_s=0.3$. The hand placement re-
 134 ward r_{hands} penalizes position and orientation errors relative to the heuristic targets. The cumulative social reward r_{social}
 135 accumulates contact-force penalties and maximum hand-placement errors over time, rewarding brief, well-located
 136 contacts while penalizing prolonged or misplaced pushes.
 137
 138

139 Training episodes include three scenario types: (1) three
 140 agents in column or line-abreast formation, with only the
 141 central agent trained; (2) single-agent omnidirectional push
 142 scenarios (as in Stage 1); and (3) a no-perturbation baseline.

143 4. Results

144 4.1. Pretraining

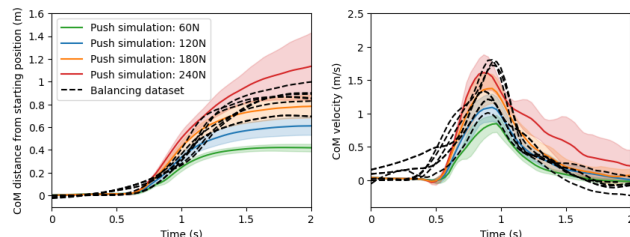


Figure 4. CoM trajectory (left) and velocity (right) for the reference (dashed) and simulated pushes at 60, 120, 180, and 240 N.

145 The policy π_{pretrain} generalizes beyond the reference
 146 dataset: Figure 4 shows CoM trajectories and velocities that
 147 match the reference data across push strengths from 60 to
 148 240 N. Stronger pushes naturally produce longer steps and,
 149 at the extremes, falls—consistent with human behavior.

150 **Ablation study.** Removing r_{balance} reduces the maximum
 151 sustainable push from ≈ 240 N to 21 N, demonstrating its
 152 critical role. Table 1 summarizes the effects of ablating r_{imit}
 153 (increased heading deviation) and r_{quality} (foot sliding and
 excess kinetic energy).

	π_{pretrain}	No r_{imit}	No r_{quality}
Heading dev.↓	5.93°	44.81°	13.19°
Foot sliding↓	22 cm	24 cm	49 cm
Kinetic energy↓	933 J	2381 J	1444 J

Table 1. Pretraining ablation over 80 trials (16 directions \times 5 force levels). Lower is better.

154

155 4.2. Adaptation

The policy π_{adapt} produces socially appropriate motions:
 156 agents raise hands to shoulder height only when a nearby
 157 agent obstructs their recovery path, and return hands to a
 158 neutral position afterward.
 159

Ablation study. Table 2 reports average hand height and
 160 transmitted impulse over 90 trials (9 directions \times 2 forma-
 161 tions \times 5 force levels). Without r_{hands} , agents keep hands
 162 raised unnecessarily after a push (final height 1.16 m vs.
 163 neutral 0.81 m). Without r_{social} , hands remain at the sides,
 164 leading to unmitigated head and torso collisions and nearly
 165 double the transmitted impulse.
 166

	π_{adapt}	No r_{hands}	No r_{social}
Final hand height↓	0.81 m	1.16 m	0.81 m
Max. hand height	0.85 m	1.16 m	0.81 m
Impulse transm.↓	40 Ns	75 Ns	74 Ns

Table 2. Adaptation ablation over 90 trials. Neutral hand height is 0.81 m. Lower is better except max. hand height (a raised hand is desired during contact).

167 4.3. Applications

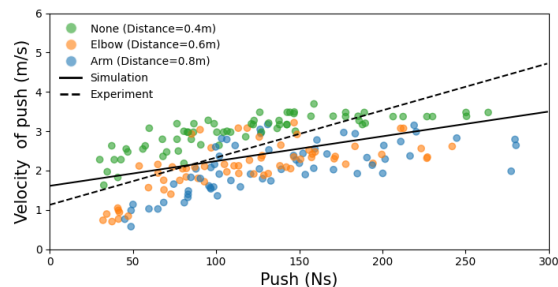


Figure 5. Push velocity (top) and propagation distance (bottom) as a function of impulse for varying interpersonal distances, consistent with empirical trends [6].

Push propagation in a line. We simulate a line of agents
 168 pushed from one end, varying interpersonal distances (0.4,
 169 0.6, 0.8 m). Figure 3 shows CoM trajectories over time;
 170 black dots mark hand-contact events, and the black line is a
 171 linear regression fit of the push-propagation wavefront. At
 172 0.4 m spacing, pushes cascade rapidly down the line with
 173 amplification; at 0.8 m, propagation stalls as agents have
 174 room to step and absorb energy. These trends are consistent
 175 with empirical studies of crowd push propagation [16, 20].
 176

177 **Dense crowd.** We simulate large crowds using π_{adapt}
178 and the hand-contact heuristic, generalizing well beyond
179 the three-agent training configurations. Figure 1 shows
180 a moving-wall force propagating through a crowd, with
181 agents adapting individually. At extreme density, small per-
182 turbations trigger rapid chaotic dispersion, reproducing am-
183 plification dynamics observed in real incidents. An addi-
184 tional ablation confirms that π_{pretrain} alone fails in multi-
185 agent settings, underscoring the necessity of the adaptation
186 stage.

187 5. Conclusion

188 We presented Push-and-Step, a physics-based framework
189 for simulating full-body humanoid agents recovering bal-
190 ance through stepping and socially aware hand-shoulder
191 contacts in dense crowds. The two-stage RL pipeline—
192 adversarial imitation learning for omnidirectional step-
193 ping, followed by AdaptNet-based finetuning for multi-
194 agent interactions—produces agents that generalize well
195 beyond their training scenarios and scale to large popula-
196 tions. Simulations reproduce empirically documented push-
197 propagation trends, demonstrating the framework’s validity
198 as a tool for studying crowd safety.

199 Current limitations include a small single-subject dataset
200 that constrains motion diversity, an engineered contact
201 heuristic that could be replaced by a learned module given
202 richer data, and a focus on static scenarios that precludes
203 dynamic crowd scenarios involving locomotion. Extending
204 the framework to navigating crowds, incorporating larger
205 multi-person interaction datasets, and replacing heuristic
206 components with learned contact policies remain promis-
207 ing directions. As a principled first step toward physically
208 grounded dense-crowd simulation, this work opens avenues
209 for applications in architecture, event planning, and crowd-
210 safety analysis.

211 References

212 [1] J. Adrian, A. Seyfried, and A. Sieben. Crowds in front of bottlenecks
213 at entrances from the perspective of physics and social psychology.
214 *Journal of the Royal Society Interface*, 17(165):20190871, 2020. 1
215 [2] Z. Aftab, T. Robert, and P.-B. Wieber. Ankle, hip and stepping
216 strategies for humanoid balance recovery with a single model pre-
217 dictive control scheme. In *IEEE-RAS International Conference on*
218 *Humanoid Robots*, pages 159–164, 2012. 1
219 [3] T. Chatagnon, S. Feldmann, J. Adrian, A.-H. Olivier, C. Pontonnier,
220 L. Hoyet, and J. Pettré. Standing balance recovery strategies of young
221 adults in a densely populated environment following external pertur-
222 bations. *Safety Science*, 177:106601, 2024. 1
223 [4] S. Coros, P. Beaudoin, and M. Van de Panne. Generalized biped
224 walking control. *ACM Transactions on Graphics*, 29(4):1–9, 2010.
225 2
226 [5] H.-T. Dang, B. Gaudou, and N. Verstaevl. A literature review of
227 dense crowd simulation. *Simulation Modelling Practice and Theory*,
228 134:102955, 2024. 1
229 [6] S. Feldmann and J. Adrian. Forward propagation of a push through
230 a row of people. *Safety science*, 164:106173, 2023. 1, 3

[7] S. Feldmann, T. Chatagnon, J. Adrian, J. Pettré, and A. Seyfried. Temporal segmentation of motion propagation in response to an external impulse. *Safety Science*, 175:106512, 2024. 2
[8] D. Helbing, I. Farkas, and T. Vicsek. Simulating dynamical features of escape panic. *Nature*, 407(6803):487–490, 2000. 1
[9] M. Kim, E. Jung, and Y. Lee. Physicsfc: Learning user-controlled skills for a physics-based football player controller. *ACM Transactions on Graphics*, 44(4):1–21, 2025. 2
[10] Y. Li, M. Lin, Z. Lin, Y. Deng, Y. Cao, and L. Yi. Learning physics-based full-body human reaching and grasping from brief walking references. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27673–27682, 2025. 2
[11] N. Pelechano, J. M. Allbeck, M. Kapadia, and N. I. Badler. *Simulating heterogeneous crowds with interactive behaviors*. CRC Press, USA, 2016. 1
[12] X. B. Peng, P. Abbeel, S. Levine, and M. Van de Panne. Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Transactions on Graphics*, 37(4):1–14, 2018. 2
[13] A. Rudenko, L. Palmieri, M. Herman, K. M. Kitani, D. M. Gavrilu, and K. O. Arras. Human motion trajectory prediction: A survey. *The International Journal of Robotics Research*, 39(8):895–935, 2020. 1
[14] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. 2
[15] A. Sieben, J. Schumann, and A. Seyfried. Collective phenomena in crowds—where pedestrian dynamics need social psychology. *PLoS one*, 12(6):e0177328, 2017. 1
[16] C. Son, D.-H. Ham, S. Jin, and T. Park. 158 deaths at halloween night: An accimap analysis of 2022 itaewon crowd crush in south korea. *Safety Science*, 184:106741, 2025. 1, 3
[17] B. J. Stephens and C. G. Atkeson. Push recovery by stepping for humanoid robots with force controlled joints. In *IEEE-RAS International Conference on Humanoid Robots*, pages 52–59, 2010. 2
[18] J. Tan, K. Liu, and G. Turk. Stable proportional-derivative controllers. *IEEE Computer Graphics and Applications*, 31(4):34–44, 2011. 2
[19] S. Tonneau, J. Pettré, and F. Multon. Using task efficient contact configurations to animate creatures in arbitrary environments. *Computers & Graphics*, 45:40–50, 2014. 1, 2
[20] W. Van Toll, T. Chatagnon, C. Braga, B. Solenthaler, and J. Pettré. SPH crowds: Agent-based crowd simulation up to extreme densities using fluid dynamics. *Computers & Graphics*, 98:306–321, 2021. 1, 3
[21] J. Won, D. Gopinath, and J. Hodgins. Control strategies for physically simulated characters performing two-player competitive sports. *ACM Transactions on Graphics*, 40(4):1–11, 2021. 2
[22] C.-C. Wu and V. Zordan. Goal-directed stepping with momentum control. In *ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, pages 113–118, 2010. 2
[23] P. Xu and I. Karamouzas. A GAN-like approach for physics-based imitation learning and interactive control. *Proceedings of the ACM on Computer Graphics and Interactive Techniques*, 4(3), 2021. 2
[24] P. Xu, Z. Wu, R. Wang, V. Sarukkai, K. Fatahalian, I. Karamouzas, V. Zordan, and C. K. Liu. Learning to ball: Composing policies for long-horizon basketball moves. *ACM Transactions on Graphics*, 44(6), 2024. 2
[25] P. Xu, K. Xie, S. Andrews, P. G. Kry, M. Neff, M. McGuire, I. Karamouzas, and V. Zordan. AdaptNet: Policy adaptation for physics-based character control. *ACM Transactions on Graphics*, 42(6), 2023. 2
[26] K. Yin, K. Loken, and M. Van de Panne. Simbicon: Simple biped locomotion control. *ACM Transactions on Graphics*, 26(3):105–es, 2007. 2
[27] V. Zordan. Angular momentum control in coordinated behaviors. In *International Conference on Motion in Games*, pages 109–120. Springer, 2010. 2