
SMPLOlympics: Sports Environments for Physically Simulated Humanoids

Anonymous Author(s)

Affiliation

<https://SMPLOlympics.github.io/>

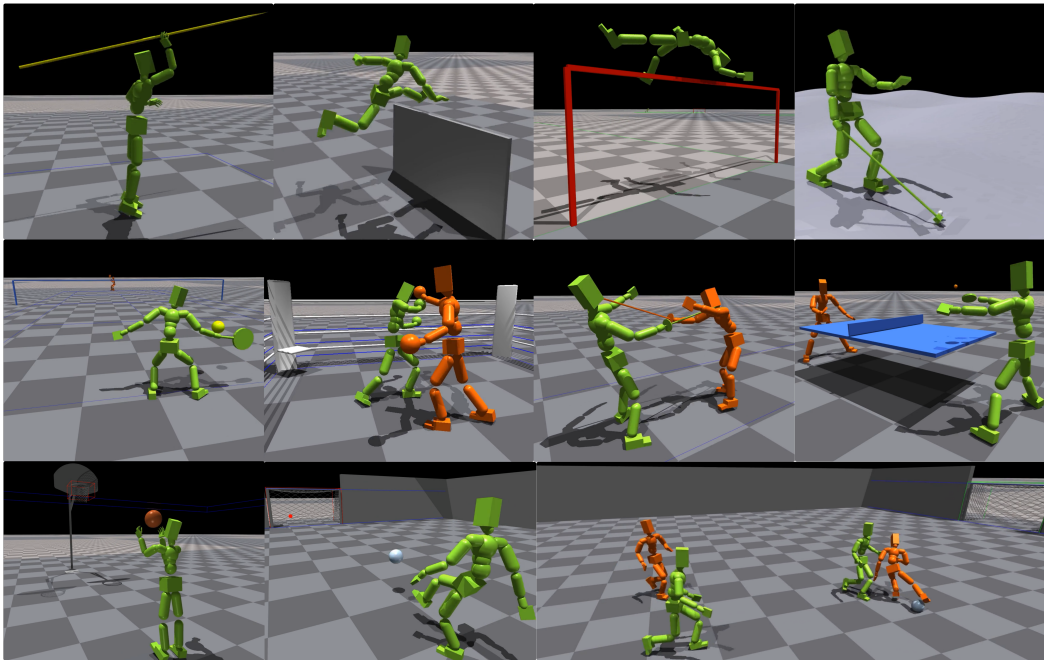


Figure 1: A collection of various sports environments for physically simulated humanoids.

Abstract

1 We present SMPLOlympics, a collection of physically simulated environments
2 that allow humanoids to compete in a variety of Olympic sports. Sports simulation
3 offers a rich and standardized testing ground for evaluating and improving the
4 capabilities of learning algorithms due to the diversity and physically demanding
5 nature of athletic activities. As humans have been competing in these sports for
6 many years, there is also a plethora of existing knowledge on the preferred strategy
7 to achieve better performance. To leverage these existing human demonstrations
8 from videos and motion capture, we design our humanoid to be compatible with
9 the widely-used SMPL and SMPL-X human models from the vision and graphics
10 community. We provide a suite of individual sports environments, including golf,
11 javelin throw, high jump, long jump, and hurdling, as well as competitive sports,
12 including both 1v1 and 2v2 games such as table tennis, tennis, fencing, boxing,
13 soccer, and basketball. Our analysis shows that combining strong motion priors
14 with simple rewards can result in human-like behavior in various sports. By
15 providing a unified sports benchmark and baseline implementation of state and
16 reward designs, we hope that SMPLOlympics can help the control and animation
17 communities achieve human-like and performant behaviors.

18 1 Introduction

19 Competitive sports, much like their role in human society, offer a standardized way of measuring
20 the performance of learning algorithms and creating emergent human behavior. While there exist
21 isolated efforts to bring individual sport into physics simulation [8, 34, 7, 33, 27], each work uses
22 a different humanoid, simulator, and learning algorithm, which prevents unified evaluation. Their
23 specially built humanoids also make it difficult to acquire compatible motion data, as retargeting
24 might be required to translate motion to each humanoid. Building a collection of simulated sports
25 environments that uses a shared humanoid embodiment and training pipeline is challenging, as it
26 requires expert knowledge in humanoid design, reinforcement learning (RL), and physics simulation.

27 These challenges have led to previous benchmarks and simulated environments [3, 25] focusing
28 mainly on locomotion tasks for humanoids. While these tasks (e.g., moving forward, getting up from
29 the ground, traversing terrains) are as benchmarks, they lack the depth and diversity needed to induce
30 a wide range of behaviors and strategies. As a result, these environments do not fully exploit the
31 potential of humanoids to discover actions and skills found in real-world human activities.

32 Another important aspect of working with simulated humanoids is the ease of obtaining human
33 demonstrations. The resemblance to the human body makes humanoids capable of performing a
34 diverse set of skills; a human can also easily judge the strategies used by humanoids. Curated human
35 motion can be used either as motion prior [17, 18, 24] or in evaluation protocols. Thus, having
36 an easy way to obtain new human motion data compatible with the humanoid, either from motion
37 capture (MoCap) or videos, is critical for simulated humanoid environments.

38 In this work, we propose SMPLOlympics, a collection of physically simulated environments for
39 a variety of Olympic sports. SMPLOlympics offers a wide range of sports scenarios that require
40 not only locomotion skills, but also manipulation, coordination, and planning. Unified under one
41 humanoid embodiment, our environments provide a rich set of challenges for developing and testing
42 embodied agents. We use humanoids compatible with the SMPL family of models, which enables
43 the direct conversion of human motion in the SMPL format to our humanoid. For tasks that require
44 articulated fingers, we use SMPL-X [16] based humanoid which has a much higher degree of
45 freedom (DOF); for tasks that do not need hands, we use SMPL [2]. As popular human models, the
46 SMPL family of models is widely adopted in the vision and graphics community, which provides
47 us with access to human pose estimation methods [32] capable of extracting coherent motion from
48 videos. The existing large-scale human motion dataset [13] in the SMPL format also helps build
49 general-purpose motion representation for humanoids [10].

50 Our sports environments support both individual and competitive sports, providing a comprehensive
51 platform for testing and benchmarking. For individual sports, we include activities such as golf,
52 javelin throw, high jump, long jump, and hurdling. Competitive sports in our suite include 1v1
53 games such as ping pong, tennis, fencing, and boxing, as well as team sports such as soccer and
54 basketball. To facilitate benchmarking, we also include tasks such as penalty kicks (for soccer) and
55 ball-target hitting (for ping-pong and tennis) that are easy to measure performance. To demonstrate
56 the importance of human demonstrations, we extract motion from videos using off-the-shelf pose
57 estimation methods, and show that using human motion data as motion prior can [18] significantly
58 improves human likeness in the resulting motion. We also test recent motion representations in
59 simulated humanoid control using hierarchical RL [10], and show that a learned motion representation
60 combined with simple rewards can lead to many versatile human-like behaviors to achieve impressive
61 sports results (*i.e.* discovering the Fosbury way for high jump).

62 In conclusion, our contributions are: (1) we propose SMPLOlympics, a collection of simulated
63 environments that allow humanoids to compete in a variety of Olympic sports; (2) we extract human
64 demonstration data from videos and show their effectiveness in helping build human-like strategies
65 in simulated sports; (3) we provide the starting state and reward designs for each sport, benchmark
66 state-of-the-art algorithms, and show that simple rewards combined with a strong motion prior can
67 lead to impressive sports feats.

68 2 Related Works

69 **Simulated Humanoid Sports.** Simulated humanoid sports can help generate animations and explore
70 optimal sports strategies. Research has focused on various individual sports within simulated
71 environments, including tennis [34], boxing [27, 36], fencing [27], basketball dribbling [7] and soccer
72 [29, 8]. These studies leverage human motion to achieve human-like behaviors, using it to acquire
73 motor skills [8, 27] or establish motion prior [34]. However, the diversity in humanoid definitions
74 across studies makes it difficult to aggregate additional human demonstration data due to the need for
75 retargetting. Furthermore, the task-specific training pipelines in these studies are hard to generalize
76 to new sports. In contrast, SMPLOlympics provides a unified benchmark employing a consistent
77 humanoid model and training pipeline across all sports. This standardization not only facilitates the
78 extension to more sports, but also simplifies the process of benchmarking learning algorithms.

79 **Simulated RL Benchmarks.** Simulated full-body humanoids provide a valuable platform for
80 studying embodied intelligence due to their close resemblance to real-world human behavior and
81 physical interactions. Current RL benchmarks [3, 25, 14] often focus on locomotion tasks such
82 as moving forward and traversing terrain. `dm_control` [25] and OpenAI [3] Gym only include
83 locomotion tasks. ASE [19] includes results for five tasks based on mocap data, which involve
84 mainly simple locomotion and sword-swinging actions. These tasks lack the complexity required
85 to fully exploit the capabilities of simulated humanoids. Sports scenarios require agile motion and
86 strategic teamwork. They are also easily interpretable and provide measurable outcomes for success.
87 A concurrent work, HumanoidBench [23] employs a commercially available humanoid robot in
88 simulation to address 27 locomotion and manipulation tasks. Unlike HumanoidBench, ours targets
89 competitive sports and uses available human demonstration data to enhance the learning of human-
90 like behaviors. This emphasis is essential, as without human demonstrations, behaviors developed in
91 benchmarks can often appear erratic, nonhuman-like, and inefficient.

92 **Humanoid Motion Representation.** Adversarial learning has proven to be a powerful method for
93 using human reference motions to enhance the naturalness of humanoid animations [18, 30, 1]. Due
94 to the high DoF in humanoids and the inherent sample inefficiency of RL training, efforts have
95 focused on developing motion primitives [6, 15, 5, 20] and motion latent spaces [4, 19, 24]. These
96 techniques aim to accelerate training and provide human-like motion priors. Notably, approaches such
97 as ASE [19], CASE [4], and CALM [24] utilize adversarial learning objectives to encourage mapping
98 between random noise and realistic motor behavior. Furthermore, methods such as ControlVAE [31],
99 NPMP [15], PhysicsVAE [28], NCP [36], and PULSE [10] leverage the motion imitation task to
100 acquire and reuse motor skills for the learning of downstream tasks. In this work, we study AMP
101 [18] and PULSE [10] as exemplary methods to provide motion priors. Our findings demonstrate
102 that a robust motion prior, combined with straightforward reward designs, can effectively induce
103 human-like behaviors in solving complex sports tasks.

104 3 Preliminaries

105 We define the full-body human pose as $\mathbf{q}_t \triangleq (\boldsymbol{\theta}_t, \mathbf{p}_t)$, consisting of 3D joint rotations $\boldsymbol{\theta}_t \in \mathbb{R}^{J \times 6}$
106 and positions $\mathbf{p}_t \in \mathbb{R}^{J \times 3}$ of all J joints on the humanoid, using the 6 DoF rotation representation
107 [35]. To define velocities $\dot{\mathbf{q}}_{1:T}$, we have $\dot{\mathbf{q}}_t \triangleq (\boldsymbol{\omega}_t, \mathbf{v}_t)$ as angular $\boldsymbol{\omega}_t \in \mathbb{R}^{J \times 3}$ and linear velocities
108 $\mathbf{v}_t \in \mathbb{R}^{J \times 3}$. If an object is involved (*e.g.* javelin, football, ping-pong ball), we define their 3D
109 trajectories $\mathbf{q}_t^{\text{obj}}$ using object position $\mathbf{p}_t^{\text{obj}}$, orientation $\boldsymbol{\theta}_t^{\text{obj}}$, linear velocity $\mathbf{v}_t^{\text{obj}}$, and angular velocity
110 $\boldsymbol{\omega}_t^{\text{obj}}$. As a notation convention, we use $\hat{\cdot}$ to denote the ground truth kinematic quantities from Motion
111 Capture (MoCap) and normal symbols without accents for values from the physics simulation.

112 **Goal-conditioned Reinforcement Learning for Humanoid Control.** We define each sport using
113 the general framework of goal-conditioned RL. Namely, a goal-conditioned policy π_{task} is trained to
114 control a simulated humanoid competing in a sports environment. The learning task is formulated
115 as a Markov Decision Process (MDP) defined by the tuple $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma \rangle$ of states, actions,

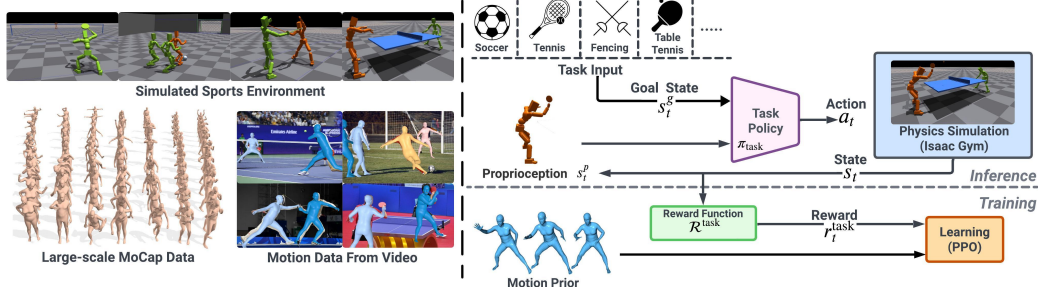


Figure 2: An overview of SMPLOlympics: we design a collection of simulated sports environments and leverage RL and human demonstrations (from videos or MoCap) as prior to tackle them.

116 transition dynamics, reward function, and discount factor. The simulation determines the state
 117 $s_t \in \mathcal{S}$ and transition dynamics T , where a policy computes the action a_t . The state s_t contains the
 118 proprioception s_t^p and the goal state s_t^g . Proprioception is defined as $s_t^p \triangleq (q_t, \dot{q}_t)$, which contains
 119 the 3D body pose q_t and velocity \dot{q}_t . We use b to indicate the boundary of the arena to which a sport
 120 is limited. All values are normalized with respect to the humanoid heading (yaw).

121 4 SMPLOlympics: sports environments For Simulated Humanoids

122 In this section, we describe the formulation of each of our sports environments, from single-person
 123 sports (Sec. 4.1) to multi-person sports (Sec. 4.2). Then, we describe our pipeline for acquiring
 124 human demonstration data from videos (Sec. 4.3). An overview can be found in Fig. 2. For each
 125 sport, we provide a preliminary reward design that serves as a baseline for future research. Due to
 126 space constraints, omitted details can be found in the supplement.

127 4.1 Single-person Sports

128 **High Jump.** In the high jump environment, the humanoid’s objective is to jump over a horizontal
 129 bar placed at a certain height without touching it. The bar is positioned following the setup of the
 130 official Olympic game. The high jump goal state $s_t^{\text{g-high_jump}} = (p_t^b, p_t^l)$ contains the positions of the
 131 bar $p_t^b \in \mathbb{R}^3$ and the landing area $p_t^l \in \mathbb{R}^3$. The reward is defined as $\mathcal{R}^{\text{high_jump}}(s_t^p, s_t^{\text{g-high_jump}}) =$
 132 $w^p r_t^p + w^h r_t^h$. The position reward r_t^p encourages the humanoid to go closer to the goal point, which
 133 is behind the high jump bar. The height reward r_t^h encourages the humanoid to jump higher. Training
 134 terminates when the humanoid is in contact with the bar, does not pass the bar, or falls to the ground
 135 before jumping. We also set up four bar heights for curriculum learning: 0.5m, 1m, 1.5m, and 2m.

136 **Long Jump.** Long jump is also set similar to the Olympic games, with a 20m runway followed
 137 by a jump area. Before the humanoid jumps, its feet should be behind the jump line. The goal
 138 state $s_t^{\text{g-long_jump}} \triangleq (p_t^s, p_t^l, p_t^g)$ includes the position of the starting point $p_t^s \in \mathbb{R}^3$, jump line
 139 $p_t^l \in \mathbb{R}^3$, and the goal $p_t^g \in \mathbb{R}^3$. The training reward is defined as $\mathcal{R}^{\text{long_jump}}(s_t^p, s_t^{\text{g-long_jump}}) \triangleq$
 140 $w^p r_t^p + w^v r_t^v + w^h r_t^h + w^l r_t^l$. The position reward r_t^p encourages the humanoid to get closer to the
 141 goal, the velocity reward r_t^v encourages larger running speed, and the height reward r_t^h encourages
 142 higher jump. Finally, r_t^l encourages jumping far.

143 **Hurdling.** In hurdling, the humanoid tries to reach a finishing line 110 meters ahead and needs to
 144 jump over 10 hurdles (each 1.067m high, placed 13.72m from the start, with subsequent hurdles
 145 spaced every 9.14m). The goal state is defined as $s_t^{\text{g-hurdling}} \triangleq (p_t^h, p_t^f)$, where $p_t^h \in \mathbb{R}^{10 \times 3}$ and
 146 $p_t^f \in \mathbb{R}^3$ includes the positions of these hurdles as well as the finish line. We define a simple reward
 147 function as $\mathcal{R}^{\text{hurdling}}(s_t^p, s_t^{\text{g-hurdling}}) = r_t^{\text{distance}}$. $\mathcal{R}^{\text{hurdling}}$ encourages the agent to run towards the finish
 148 line and clear each hurdle. Additionally, we employ a curriculum for hurdling, where the height of
 149 each hurdle is randomly sampled between 0 and 1.167 meters for each episode.

150 **Golf.** For golf, the humanoid’s right hand is replaced with a golf club measuring 1.14 meters. The
 151 driver of the golf club is simulated as a small box (0.05m \times 0.025m \times 0.02m). We incorporate a

152 randomly generated terrain in the golf environment, designed to mimic real-world grasslands with
 153 wave-like features and an amplitude of 0.5 meters. The objective for the humanoid is to hit the ball
 154 towards a randomly sampled target position. The goal state $\mathbf{s}_t^{\text{g-golf}} \triangleq (\mathbf{p}_t^b, \mathbf{p}_t^c, \mathbf{p}_t^g, \mathbf{o}_t)$ includes the ball
 155 position $\mathbf{p}_t^b \in \mathbb{R}^3$, club $\mathbf{c}_t^b \in \mathbb{R}^3$, goal position $\mathbf{p}_t^g \in \mathbb{R}^3$, and terrain height map $\mathbf{o}_t \in \mathbb{R}^{32 \times 32}$. The
 156 reward is defined as $\mathcal{R}^{\text{golf}}(\mathbf{s}_t^p, \mathbf{s}_t^{\text{g-golf}}) = w^p r_t^p + w^c r_t^c + w^g r_t^g + w^{\text{pred}} r_t^{\text{pred}}$, where the r_t^p encourages
 157 the ball to move forward, r_t^c encourages swinging the golf club to hit the ball, and r_t^g encourages the
 158 ball to reach the goal. In addition, we predict the ball’s trajectory and provide a dense reward r_t^{pred}
 159 based on the distance between the predicted landing point and the goal.

160 **Javelin.** For javelin throw, we use SMPL-X humanoid with articulated fingers. The goal state is
 161 defined as $\mathbf{s}_t^{\text{g-javelin}} \triangleq (\mathbf{q}_t^{\text{obj}}, \mathbf{p}_t^r, \mathbf{p}_t^h)$, where $\mathbf{q}_t^{\text{obj}} \in \mathbb{R}^{13}$, includes the position, orientation, linear, and
 162 angular velocity of the javelin. \mathbf{p}_t^r and \mathbf{p}_t^h are the positions of the root and right hand. The reward
 163 is defined as $\mathcal{R}^{\text{javelin}}(\mathbf{s}_t^p, \mathbf{s}_t^{\text{g-javelin}}) \triangleq w^{\text{grab}} r_t^{\text{grab}} + w^{\text{js}} r_t^{\text{js}} + w^{\text{goal}} r_t^{\text{goal}} + w^s r_t^s$. The grab reward r_t^{grab}
 164 encourages the right hand to grab the javelin. The javelin stability reward r_t^{js} minimizes the javelin’s
 165 self-rotation. The goal reward r_t^{goal} encourages the humanoid to throw the javelin further. The stability
 166 reward r_t^s is to avoid large movements of the body.

167 4.2 Multi-person Sports

168 **Tennis.** For tennis, each humanoid’s right hand is replaced as an oval racket. We use the same
 169 measurement as a real tennis court and ball. We design two tasks: a single-player task where the
 170 humanoid trains to hit balls launched randomly, and a 1v1 mode where the humanoid plays against
 171 another humanoid. For both tasks, the goal state is defined as $\mathbf{s}_t^{\text{g-tennis}} \triangleq (\mathbf{p}_t^{\text{ball}}, \mathbf{v}_t^{\text{ball}}, \mathbf{p}_t^{\text{racket}}, \mathbf{p}_t^{\text{tar}}$,
 172 where $\mathbf{p}_t^{\text{ball}} \in \mathbb{R}^3$, $\mathbf{v}_t^{\text{ball}} \in \mathbb{R}^3$, $\mathbf{p}_t^{\text{racket}} \in \mathbb{R}^3$ and $\mathbf{p}_t^{\text{tar}} \in \mathbb{R}^3$, which includes the position and velocity of
 173 the ball, position of the racket and position of the target. The reward function for tennis is defined
 174 as $\mathcal{R}^{\text{tennis}}(\mathbf{s}_t^p, \mathbf{s}_t^{\text{g-tennis}}) = w_p r_t^{\text{racket}} + w_b r_t^{\text{ball}}$. The racket reward r_t^{racket} encourages the racket to reach
 175 the ball, and the ball reward r_t^{ball} aims to successfully hit the ball into the opponent’s court, as close
 176 to the target as possible. For the single-player task, we shoot a ball from the opposite side from a
 177 random position and trajectory, simulating a ball hit by the opponent. The target $\mathbf{p}_t^{\text{tar}}$ is also randomly
 178 sampled. For the 1v1 scenario, we can either train models from scratch or initialize two identical
 179 single-player models as opponents, which can play back and forth.

180 **Table Tennis.** For table tennis, each humanoid is equipped with a circular paddle (replacing the right
 181 hand) and play on a standard table. Similar to tennis, we have the single-player task and the 1v1 task.
 182 Similarly, the goal state is defined as $\mathbf{s}_t^{\text{g-table-tennis}} \triangleq (\mathbf{p}_t^{\text{ball}}, \mathbf{v}_t^{\text{ball}}, \mathbf{p}_t^{\text{racket}}, \mathbf{p}_t^{\text{tar}})$. The reward function for
 183 table tennis is defined as $\mathcal{R}^{\text{table-tennis}}(\mathbf{s}_t^p, \mathbf{s}_t^{\text{g-table-tennis}}) = w_p r_t^{\text{racket}} + w_b r_t^{\text{ball}}$. The paddle reward r_t^{racket}
 184 is the same as the tennis while we modify the r_t^{ball} slightly to encourage more hits for table tennis.

185 **Fencing.** For 1v1 fencing, each humanoid is equipped with a sword (replacing the right hand)
 186 and plays on a standard fencing field. The goal state is defined as $\mathbf{s}_t^{\text{g-fencing}} \triangleq (\mathbf{p}_t^{\text{opp}}, \mathbf{v}_t^{\text{opp}}, \mathbf{p}_t^{\text{sword}} -$
 187 $\mathbf{p}_t^{\text{opp-target}}, \|\mathbf{c}_t\|_2^2, \|\mathbf{c}_t^{\text{opp}}\|_2^2, \mathbf{b})$, which contains the opponent’s position body $\mathbf{p}_t^{\text{opp}} \in \mathbb{R}^{24 \times 3}$, linear
 188 velocity $\mathbf{v}_t^{\text{opp}} \in \mathbb{R}^{24 \times 3}$, the difference between target body position $\mathbf{p}_t^{\text{opp-target}} \in \mathbb{R}^{5 \times 3}$ on the opponent
 189 and agent’s sword tip position $\mathbf{p}_t^{\text{sword}}$, normalized contract forces on the agent itself $\|\mathbf{c}_t\|_2^2 \in \mathbb{R}^{24 \times 3}$
 190 and its opponent $\|\mathbf{c}_t^{\text{opp}}\|_2^2 \in \mathbb{R}^{24 \times 3}$, as well as the bounding box $\mathbf{b} \in \mathbb{R}^4$. To train the fencing agent, we
 191 define the fencing reward function as $\mathcal{R}^{\text{fencing}}(\mathbf{s}_t^p, \mathbf{s}_t^{\text{g-fencing}}) = w_f r_t^{\text{facing}} + w_v r_t^{\text{vel}} + w_s r_t^{\text{strike}} + w_p r_t^{\text{point}}$.
 192 The facing r_t^{facing} and velocity reward r_t^{vel} encourage the agent to face and move toward the opponent.
 193 The strike reward r_t^{strike} encourages the agent’s sword tip to get close to the target, while r_t^{point} is the
 194 reward for getting in contact with the target. We use the pelvis, head, spine, chest, and torso as the
 195 target bodies. The episode terminates if either of the humanoids falls or steps out of bounds.

196 **Boxing.** For boxing, we simulate two humanoids with sphere hands in a bounded arena. The goal
 197 state is similar to fencing: $\mathbf{s}_t^{\text{g-boxing}} \triangleq (\mathbf{p}_t^{\text{opp}}, \mathbf{v}_t^{\text{opp}}, \mathbf{p}_t^{\text{hand}} - \mathbf{p}_t^{\text{opp-target}}, \|\mathbf{c}_t\|_2^2, \|\mathbf{c}_t^{\text{opp}}\|_2^2)$ but without the
 198 bounding box information. The reward function and target body parts are also the same as fencing,
 199 though replacing the sword tip to the hands.

200 **Soccer.** The soccer environment includes one or more humanoids, a ball, two goal posts, and the field
 201 boundaries. The field measures 32m \times 20m. We support three tasks: penalty kicks, 1v1, and 2v2.

202 For penalty kicks, the humanoid is positioned 13 meters from the goal line, with the ball placed
 203 at a fixed spot 12 meters directly in front of the goal center. The objective is to kick the ball
 204 toward a randomly sampled target within the goal post. To achieve this, the controller is provided
 205 $s_t^{\text{g-kick}} \triangleq (\mathbf{p}_t^{\text{ball}}, \dot{\mathbf{q}}_t^{\text{ball}}, \mathbf{p}_t^{\text{goal-post}}, \mathbf{p}_t^{\text{goal-target}})$, where $\mathbf{p}_t^{\text{ball}} \in \mathbb{R}^3$ is the ball position, $\dot{\mathbf{q}}_t^{\text{ball}} \in \mathbb{R}^3$ is the
 206 velocity and angular velocity, $\mathbf{p}_t^{\text{goal-post}} \in \mathbb{R}^4$ is the bounding box of the goal, and $\mathbf{p}_t^{\text{goal-target}} \in \mathbb{R}^3$ is
 207 the target location within the goal post. The reward is $\mathcal{R}^{\text{soccer-kick}}(s_t^{\text{p}}, s_t^{\text{g-kick}}) \triangleq w^{\text{p}2\text{b}} r^{\text{p}2\text{b}} + w^{\text{b}2\text{g}} r^{\text{b}2\text{g}} +$
 208 $w^{\text{bv}2\text{g}} r^{\text{bv}2\text{g}} + w^{\text{b}2\text{t}} r^{\text{b}2\text{t}} - c_t^{\text{no-dribble}}$. Various rewards are designed to guide the character towards a
 209 run-and-kick motion. The player-to-ball ($r^{\text{p}2\text{b}}$) reward motivates the character to move towards the
 210 ball. The ball-to-goal reward ($r^{\text{b}2\text{g}}$) reduces the distance between the ball and the target. The ball-
 211 velocity-to-goal ($r^{\text{bv}2\text{g}}$) encourages a higher velocity of the ball toward the target. The ball-to-target
 212 ($r^{\text{b}2\text{t}}$) reward encourages a smaller distance between the target and the predicted landing spot of the
 213 ball based on its current position and velocity. Finally, a negative reward ($c_t^{\text{no-dribble}}$) is applied if the
 214 character passes the spawn position of the ball, which discourages dribbling and encourages kicking.

215 Beyond penalty kicks, we explore team-play dynamics, including 1v1 and 2v2. The controller is
 216 provided with a state defined as $s_t^{\text{g-soccer}} \triangleq (\mathbf{p}_t^{\text{ball}}, \dot{\mathbf{q}}_t^{\text{ball}}, \mathbf{p}_t^{\text{goal-post}}, \mathbf{p}_t^{\text{ally-root}}, \mathbf{p}_t^{\text{opp-root}})$, where $\mathbf{p}_t^{\text{ally-root}} \in$
 217 \mathbb{R}^3 and $\mathbf{p}_t^{\text{opp-root}} \in \mathbb{R}^3$ are the root positions of the ally and opponents (1 or 2). The controller is then
 218 trained using the following reward $\mathcal{R}^{\text{soccer-match}}(s_t^{\text{p}}, s_t^{\text{g-soccer}}) \triangleq w^{\text{p}2\text{b}} r^{\text{p}2\text{b}} + w^{\text{b}2\text{g}} r^{\text{b}2\text{g}} + w^{\text{bv}2\text{g}} r^{\text{bv}2\text{g}} +$
 219 $w^{\text{point}} r^{\text{point}}$, where $r^{\text{p}2\text{b}}$, $r^{\text{b}2\text{g}}$ and $r^{\text{bv}2\text{g}}$ are the same as in penalty kick. $r^{\text{b}2\text{g}}$ and $r^{\text{bv}2\text{g}}$ are zeroed out
 220 when the distance to the ball is greater than 0.5m. r^{point} , the scoring a goal, provides a one-time bonus
 221 and or penalty for goals. Notice that this is a rudimentary reward design compared to prior art [8] and
 222 serves as a starting point for further development.

223 **Basketball.** Our basketball environment is set up similarly to the soccer environment except for using
 224 the SMPL-X humanoid. The court measures 29m \times 15m, with a 3m high hoop. We also introduce
 225 the task of free-throw, where the humanoid begins at a distance of 4.5 meters from the hoop with the
 226 ball initially positioned close to its hands. The objective is to successfully throw the basketball into
 227 the hoop. The goal state for this task is defined similarly to that of the soccer penalty kicks, with the
 228 distinction being the prohibition of foot-to-ball contact to maintain basketball rules.

229 **Competitive Self-play.** In competitive sports environments, we implement a basic adversarial self-
 230 play mechanism where two policies, initialized randomly, compete against each other to optimize
 231 their rewards. We adopt an alternating optimization strategy from [27], where one policy is frozen
 232 while the other is trained. This encourages each policy to develop offensive and defensive strategies,
 233 contributing to more competitive behavior, as observed in boxing and fencing ([supplement site](#)).

234 4.3 Acquiring Human Demonstration From Videos

235 We utilize TRAM [26] for 3D motion reconstruction from Internet videos, providing robust global
 236 trajectory and pose estimation under dynamic camera movements, commonly found in sports broad-
 237 casting. Specifically, TRAM estimates SMPL parameters [9] which include global root translation,
 238 orientation, body poses, and shape parameters. We further apply PHC [11], a physics-based motion
 239 tracker, to imitate these estimated motions, ensuring physical plausibility. We find these corrected
 240 motions are significantly more effective as positive samples for adversarial learning compared to raw
 241 estimated results. More details and ablation are provided in the supplementary materials.

242 5 Experiments

243 **Implementation Details.** Simulation is conducted in Isaac Gym [14], where the policy runs at 30
 244 Hz and the simulation at 60 Hz. All task policies utilize three-layer MLPs with units [2048, 1024,
 245 512]. The SMPL humanoid models adhere to the SMPL kinematic structure, featuring 24 joints,
 246 23 of which are actuated, yielding an action space of \mathcal{R}^{69} . The SMPL-X humanoid has 52 joints,

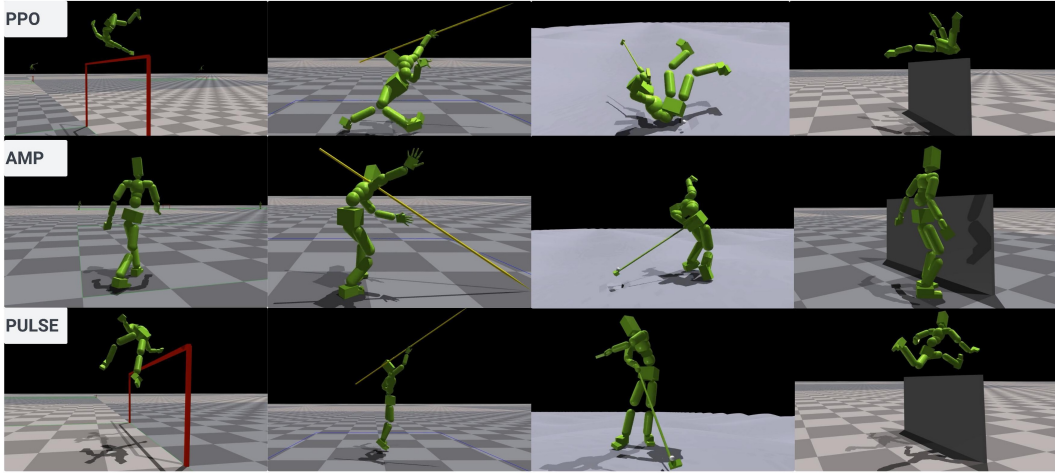


Figure 3: Qualitative results for high jump, javelin, golf, and hurdling. PPO and AMP try to solve the task using inhuman behavior, while PULSE can discover human-like behavior.

247 51 actuated, including 21 body joints and hands, resulting in an action space of \mathcal{R}^{153} . Body parts
 248 on our humanoid consist of primitives such as capsules and blocks. All models can be trained on a
 249 single Nvidia RTX 3090 GPU in 1-3 days. We limit all joint actuation forces to 500 Nm. For more
 250 implementation details, please refer to the supplement.

251 **Baselines.** We benchmark our simulated sports using some of the state-of-the-art simulated humanoid
 252 control methods. While not a comprehensive list, it provides a baseline for the challenging environ-
 253 ments. Each task is trained using PPO [22], AMP [18], PULSE [10], and a combination of PULSE
 254 and AMP. AMP use a discriminator with the policy to provide an adversarial reward, using human
 255 demonstration data to deliver a “style” reward that reflects the human-likeness of humanoid motion.
 256 Both task and discriminator rewards are equally weighted at 0.5. PULSE extracts a 32-dimensional
 257 universal motion representation from AMASS data, surpassing previous methods [24, 19] in coverage
 258 of motor skills and applicability to downstream tasks. Compared to AMP, PULSE uses hierarchical
 259 RL and offers a learned action space that accelerates training and provides human-like motion prior
 260 (instead of a discriminative reward). PULSE and AMP can be combined effectively, where PULSE
 261 provides the action space and AMP provides task-specific style reward.

262 **Metrics.** We provide quantitative evaluations for tasks with easily measurable metrics such as high
 263 jump, long jump, hurdling, javelin, golf, single-player tennis, table tennis, penalty kicks, and free
 264 throws. These metrics are detailed in the supplementary materials, where we also present qualitative
 265 assessments for tasks that are more challenging to quantify, such as boxing, fencing, and team soccer.
 266 Specifically, success rate (Suc Rate) determines whether an agent completes a sport according to set
 267 rules. Average distance (Avg Dis) indicates the extent an agent or object travels. For sports involving
 268 ball hits, such as tennis and table tennis, we record the average number of successful ball strikes (Avg
 269 Hits). Error distance (Error Dis) measures the distance between the intended target and the actual
 270 landing spot, applicable in sports like golf, tennis, and penalty kicks. Additionally, the hit rate in golf
 271 quantifies the success of striking the ball with the club. Evaluations are performed on 1000 trials.

272 5.1 Benchmarking Popular Simulated Humanoid Algorithms

273 In this section, we evaluate the performance of various control methods across our sports environments.
 274 We provide qualitative results in Fig. 3 and Fig. 4, and training curves in Fig. 5. To view extensive
 275 qualitative results, including human-like soccer kick, boxing, high jump, *etc.*, please see [supplement](#).

276 **Track & Field Sports (Without Video Data).** We first evaluate track and field sports, including
 277 long jump, high jump, hurdling, and javelin throwing. For these sports, SOTA pose estimation
 278 methods fail to estimate coherent motion and global root trajectory as players and cameras are both
 279 fast-moving. Thus, we utilize a subset of the AMASS dataset containing locomotion data [21] as

Table 1: Evaluation on Long Jump, High Jump, Hurdling and Javelin. World records are in parentheses.

Method	Long Jump (8.95m)			High Jump (2.45m)			Hurdling (12.8s)			Javelin (104.8m)	
	Suc Rate ↑	Avg Dis ↑	Suc Rate (1m) ↑	Height (1m) ↑	Suc Rate (1.5m) ↑	Height (1.5m) ↑	Suc Rate ↑	Avg Dis ↑	Time ↓	Suc Rate ↑	Avg Dis ↑
PPO [22]	53.6%	19.42	100%	4.08	100%	4.11	57.6%	108.9	11.22	100%	44.5
AMP [18]	0%	-	0%	-	0%	-	0%	13.24	-	0.31%	2.03
PULSE [10]	100%	5.105	100%	2.01	100%	1.98	100%	122.1	17.76	100%	9.63

Table 2: Evaluation on Golf, Tennis, Table Tennis, Penalty Kick and Free Throw

Method	Golf		Tennis		Table Tennis		Penalty Kick		Free Throw
	Hit Rate ↑	Error Dis ↓	Avg Hits ↑	Error Dis ↓	Avg Hits ↑	Error Dis ↓	Suc Rate ↑	Error Dis ↓	Suc Rate ↑
PPO [22]	0%	-	2.76	1.92	1.01	0.06	0.0%	-	0.0%
AMP [18]	100%	1.43	3.95	5.30	1.10	0.13	0.0%	-	0.0%
PULSE [10]	99.9%	1.29	2.48	3.50	0.74	0.19	76.6%	0.25	87.5%
PULSE [10] + AMP [18]	99.9%	2.18	2.62	3.64	1.83	0.23	27.5%	0.27	30.6%

reference motions. Since PULSE is pretrained on AMASS, we exclude PULSE + AMP from these tests. Table 1 summarizes the quantitative results of different methods. In long jump, AMP fails entirely, often walking slowly to the jump line without a forward leap. This failure occurs because the policy prioritizes discriminator rewards over task completion. If the task is too hard, the policy will use simple motion (such as standing still) to maximize the discriminator reward instead of trying to complete the task. In contrast, PPO, while capable of jumping great distances, exhibits unnatural motions. PULSE successfully executes jumps with human-like motion, but lacks the specialized skills for top-tier records due to the absence of corresponding motion data in AMASS. The high jump displays similar patterns: PPO achieves impressive heights but with unnatural movements while AMP struggles to reconcile adversarial and task rewards. Surprisingly, as shown in Figure 3, PULSE successfully adopts a Fosbury flop approach without specific rewards to encourage this technique, likely leveraging breakdance skills. For hurdling, AMP completely fails, stopping before the first hurdle. PPO bounces energetically over each obstacle as shown in Figure 3, but sometimes falls and fails to complete the race, with an average success rate of just over 50% and an average distance of less than 110m. PULSE facilitates natural clearance of hurdles, and completes races in 17.76 seconds, a competitive time compared to human standards. Javelin throwing poses similar challenges: PPO uses inhuman strategies, AMP struggles with balancing rewards, and PULSE adopts human-like strategies but lacks specific skills for record-setting performance.

Sports With Video Data. For sports including golf, tennis, table tennis, and soccer penalty kick, we utilize processed motion from videos as demonstrations for AMP and PULSE+AMP. The results are reported in Table 2 and Fig. 4. In tennis, AMP demonstrates superior performance in terms of average hits; however, returned balls often land far from the intended targets. This is because prolonged rallies increase discriminator rewards, leading AMP to ignore task rewards. Notably, AMP exhibits inhuman motions at the moment of ball contact and reverts to natural movements when preparing for the next hit as shown in Fig. 4. This behavior underscores a reward conflict between balancing task and discriminator rewards. PPO plays tennis in an unnatural way, while PULSE and PULSE + AMP show similar performance. In table tennis, PPO achieves impressive error distances, but struggles with consistency and often fails to return second shots. We observe video data proves *particularly beneficial for table tennis*. PULSE+AMP records significantly higher hit averages with reasonable error distances. Table tennis requires quick reactions within a short time, which the pre-trained PULSE model supports by providing necessary motor skills, enhanced by video data that guide the learning of proper stroke techniques. For golf, penalty kicks, and free throws, the “initiating contact with an object” part makes them challenging. Here, only PULSE and PULSE+AMP manage to solve the three tasks effectively, leveraging PULSE’s latent space for effective exploration. The design of these tasks often results in a sparse exploration phase where triggering penalty rewards, such as $c_t^{\text{no-dribble}}$ for moving past the ball’s initial position. The AMP reward also negatively affects training penalty kick, as the human demonstration contains other soccer motions such as running and dribbling, and the policy finds them easier to learn and exploit.

Curriculum learning. We find curriculum learning is an essential component in achieving better results for some tasks. In Table 3, we study variants of high jump and hurdling task with and without

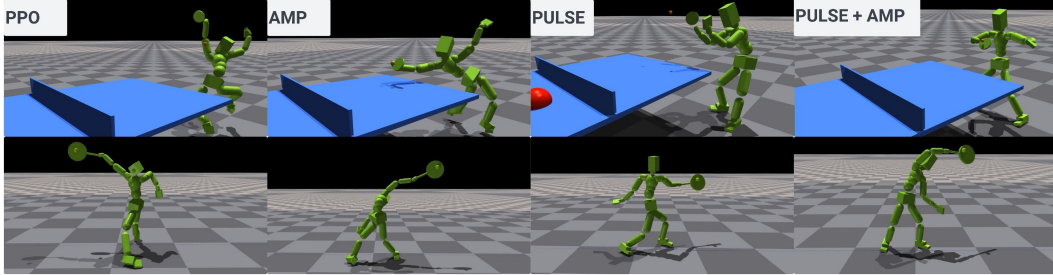


Figure 4: Qualitative results for table tennis and tennis. PPO and AMP result in inhuman behavior; PULSE can use human-like movement but PULSE + AMP result in behavior specific to the sport.

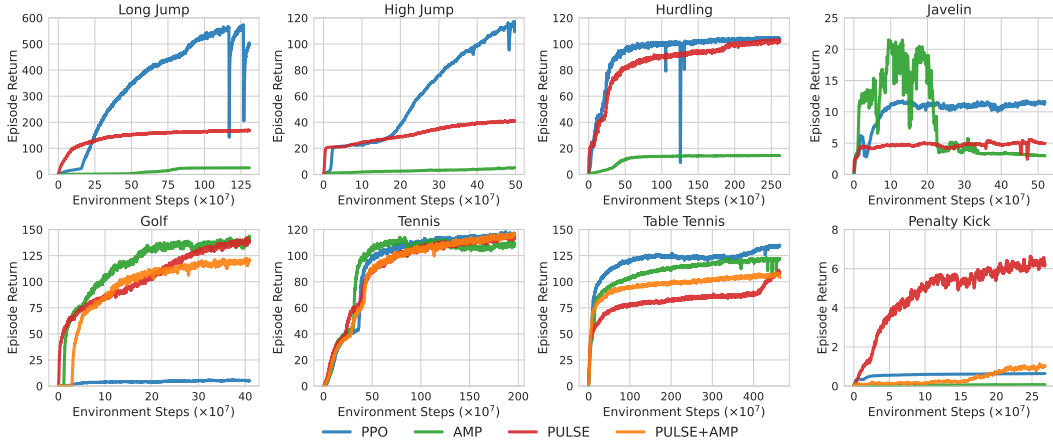


Figure 5: Learning curves on various tasks.

320 the curriculum using PULSE. We can see that
 321 without curriculum, high jump and hurdling
 322 both fail to solve the task. This is due to the
 323 policy not being able to obtain any reward fac-
 324 ing challenging heights of bars and hurdles and the policy gets stuck in the local minima.

Table 3: Evaluation on curriculum learning.

Method	High Jump		Hurdling		
	Suc Rate (1m)	Suc Rate (1.5m)	Suc Rate	Avg Dis	Time
w/o curriculum	100%	0%	0%	13.65	-
w/ curriculum	100%	100%	100%	122.1	17.76

325 6 Limitations, Conclusion and Future Work

326 **Limitations** . While SMPLOlympics provides a large collection of simulated sports environments, it
 327 is far from being comprehensive. Certain sports are omitted due to simulation constraints (e.g., swim-
 328 ming, shooting, ice hockey, cycling) or their inherent complexity (e.g., 11-a-side soccer, equestrian
 329 events). Nevertheless, our framework is highly adaptable, allowing easy incorporation of additional
 330 sports like climbing, rugby, wrestling *etc*. Our initial design of rewards, though able to achieve
 331 sensible results, is also far from optimal. For competitive sports such as 2v2 soccer and basketball,
 332 our results also fall short of SOTA [8] which employs much more complex systems.

333 **Conclusion and Future Work**. We introduce SMPLOlympics, a collection of sports environments
 334 for simulated humanoids. We provide carefully designed state and reward, and benchmark humanoid
 335 control algorithms and motion priors. We find that by combining simple reward design and powerful
 336 human motion prior, one can achieve human-like behavior for solving various challenging sports.
 337 Our humanoid’s compatibility with the SMPL family of models also provides an easy way to obtain
 338 additional data from video for training, which we demonstrate to be helpful in training some sports.
 339 These well-defined simulation environments could also serve as valuable platforms for frontier models
 340 [12] to gain physical understanding. We believe that SMPLOlympics provides a valuable starting
 341 point for the community to further explore physically simulated humanoids.

References

- [1] Jinseok Bae, Jungdam Won, Donggeun Lim, Cheol-Hui Min, and Young Min Kim. Pmp: Learning to physically interact with environments using part-wise motion priors. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–10, 2023.
- [2] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J Black. Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. *Lect. Notes Comput. Sci.*, 9909 LNCS:561–578, 2016. ISSN 0302-9743,1611-3349.
- [3] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym, 2016.
- [4] Zhiyang Dou, Xuelin Chen, Qingnan Fan, Taku Komura, and Wenping Wang. C· ase: Learning conditional adversarial skill embeddings for physics-based characters. In *SIGGRAPH Asia 2023 Conference Papers*, pages 1–11, 2023.
- [5] Tuomas Haarnoja, Kristian Hartikainen, Pieter Abbeel, and Sergey Levine. Latent space policies for hierarchical reinforcement learning. *arXiv preprint arXiv:1804.02808*, 2018.
- [6] Leonard Hasenclever, Fabio Pardo, Raia Hadsell, Nicolas Heess, and Josh Merel. CoMic: Complementary task learning & mimicry for reusable skills. In Hal Daumé Iii and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 4105–4115. PMLR, 2020.
- [7] Libin Liu and Jessica Hodgins. Learning basketball dribbling skills using trajectory optimization and deep reinforcement learning. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018.
- [8] Siqi Liu, Guy Lever, Zhe Wang, Josh Merel, S M Ali Eslami, Daniel Hennes, Wojciech M Czarnecki, Yuval Tassa, Shayegan Omidshafiei, Abbas Abdolmaleki, Noah Y Siegel, Leonard Hasenclever, Luke Marris, Saran Tunyasuvunakool, H Francis Song, Markus Wulfmeier, Paul Muller, Tuomas Haarnoja, Brendan D Tracey, Karl Tuyls, Thore Graepel, and Nicolas Heess. From motor control to team play in simulated humanoid football. *arXiv preprint arXiv:2105.12196*, 2021.
- [9] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. *ACM Trans. Graph.*, 34, 2015. ISSN 0730-0301,1557-7368.
- [10] Zhengyi Luo, Jinkun Cao, Josh Merel, Alexander Winkler, Jing Huang, Kris Kitani, and Weipeng Xu. Universal humanoid motion representations for physics-based control. *arXiv preprint arXiv:2310.04582*, 2023.
- [11] Zhengyi Luo, Jinkun Cao, Alexander W. Winkler, Kris Kitani, and Weipeng Xu. Perpetual humanoid control for real-time simulated avatars. In *International Conference on Computer Vision (ICCV)*, 2023.
- [12] Yecheng Jason Ma, William Liang, Guanzhi Wang, De-An Huang, Osbert Bastani, Dinesh Jayaraman, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Eureka: Human-level reward design via coding large language models. *arXiv preprint arXiv:2310.12931*, 2023.
- [13] Naureen Mahmood, Nima Ghorbani, Nikolaus F Troje, Gerard Pons-Moll, and Michael J Black. Amass: Archive of motion capture as surface shapes. *Proceedings of the IEEE International Conference on Computer Vision*, 2019-Octob:5441–5450, 2019. ISSN 1550-5499.
- [14] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, and Gavriel State. Isaac gym: High performance gpu-based physics simulation for robot learning. *arXiv preprint arXiv:2108.10470*, 2021.

- 387 [15] Josh Merel, Leonard Hasenclever, Alexandre Galashov, Arun Ahuja, Vu Pham, Greg Wayne,
388 Yee Whye Teh, and Nicolas Heess. Neural probabilistic motor primitives for humanoid control,
389 2018. ISSN 2331-8422.
- 390 [16] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman,
391 Dimitrios Tzionas, and Michael J. Black. Expressive body capture: 3d hands, face, and body
392 from a single image. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition*
393 *(CVPR)*, 2019.
- 394 [17] Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. Deepmimic. *ACM*
395 *Trans. Graph.*, 37:1–14, 2018. ISSN 0730-0301.
- 396 [18] Xue Bin Peng, Ze Ma, Pieter Abbeel, Sergey Levine, and Angjoo Kanazawa. Amp: Adversarial
397 motion priors for stylized physics-based character control. *ACM Trans. Graph.*, pages 1–20,
398 2021.
- 399 [19] Xue Bin Peng, Yunrong Guo, Lina Halper, Sergey Levine, and Sanja Fidler. Ase: Large-
400 scale reusable adversarial skill embeddings for physically simulated characters. *arXiv preprint*
401 *arXiv:2205.01906*, 2022.
- 402 [20] Dushyant Rao, Fereshteh Sadeghi, Leonard Hasenclever, Markus Wulfmeier, Martina Zambelli,
403 Giulia Vezzani, Dhruva Tirumala, Yusuf Aytar, Josh Merel, Nicolas Heess, and Raia Hadsell.
404 Learning transferable motor skills with hierarchical latent mixture policies. *arXiv preprint*
405 *arXiv:2112.05062*, 2021.
- 406 [21] Davis Rempe, Zhengyi Luo, Xue Bin Peng, Ye Yuan, Kris Kitani, Karsten Kreis, Sanja Fidler,
407 and Or Litany. Trace and pace: Controllable pedestrian animation via guided trajectory diffusion.
408 *arXiv preprint arXiv:2304.01893*, 2023.
- 409 [22] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal pol-
410 icy optimization algorithms, 2017. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:28695052)
411 [28695052](https://api.semanticscholar.org/CorpusID:28695052).
- 412 [23] Carmelo Sferrazza, Dun-Ming Huang, Xingyu Lin, Youngwoon Lee, and Pieter Abbeel. Hu-
413 manoidbench: Simulated humanoid benchmark for whole-body locomotion and manipulation.
414 *arXiv preprint arXiv:2403.10506*, 2024.
- 415 [24] Chen Tessler, Israel Yoni Kasten, Israel Yunrong Guo, and Canada Nvidia. Calm: Conditional
416 adversarial latent models for directable virtual characters.
- 417 [25] Saran Tunyasuvunakool, Alistair Muldal, Yotam Doron, Siqi Liu, Steven Bohez, Josh Merel,
418 Tom Erez, Timothy Lillicrap, Nicolas Heess, and Yuval Tassa. dm_control: Software and tasks
419 for continuous control. *Software Impacts*, 6:100022, 2020.
- 420 [26] Yufu Wang, Ziyun Wang, Lingjie Liu, and Kostas Daniilidis. Tram: Global trajectory and
421 motion of 3d humans from in-the-wild videos. *arXiv preprint arXiv:2403.17346*, 2024.
- 422 [27] Jungdam Won, Deepak Gopinath, and Jessica Hodgins. Control strategies for physically
423 simulated characters performing two-player competitive sports. *ACM Trans. Graph.*, 40:1–11,
424 2021. ISSN 0730-0301.
- 425 [28] Jungdam Won, Deepak Gopinath, and Jessica Hodgins. Physics-based character controllers
426 using conditional vaes. *ACM Trans. Graph.*, 41:1–12, 2022. ISSN 0730-0301.
- 427 [29] Zhaoming Xie, Sebastian Starke, Hung Yu Ling, and Michiel van de Panne. Learning soccer
428 juggling skills with layer-wise mixture-of-experts. In *ACM SIGGRAPH 2022 Conference*
429 *Proceedings*, pages 1–9, 2022.
- 430 [30] Pei Xu, Xiumin Shang, Victor Zordan, and Ioannis Karamouzas. Composite motion learning
431 with task control. *ACM Transactions on Graphics (TOG)*, 42(4):1–16, 2023.

- 432 [31] Heyuan Yao, Zhenhua Song, Baoquan Chen, and Libin Liu. Controlvae: Model-based learning
433 of generative controllers for physics-based characters. *arXiv preprint arXiv:2210.06063*, 2022.
- 434 [32] Vickie Ye, Georgios Pavlakos, Jitendra Malik, and Angjoo Kanazawa. Decoupling human
435 and camera motion from videos in the wild. In *Proceedings of the IEEE/CVF conference on*
436 *computer vision and pattern recognition*, pages 21222–21232, 2023.
- 437 [33] Zhiqi Yin, Zeshi Yang, Michiel Van De Panne, and KangKang Yin. Discovering diverse athletic
438 jumping strategies. *ACM Transactions on Graphics (TOG)*, 40(4):1–17, 2021.
- 439 [34] Haotian Zhang, Ye Yuan, Viktor Makoviychuk, Yunrong Guo, Sanja Fidler, Xue Bin Peng, and
440 Kayvon Fatahalian. Learning physically simulated tennis skills from broadcast videos. *ACM*
441 *Trans. Graph.*, 42:1–14, 2023. ISSN 0730-0301,1557-7368.
- 442 [35] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation
443 representations in neural networks. *Proceedings of the IEEE Computer Society Conference on*
444 *Computer Vision and Pattern Recognition*, 2019-June:5738–5746, 2019. ISSN 1063-6919.
- 445 [36] Qingxu Zhu, He Zhang, Mengting Lan, and Lei Han. Neural categorical priors for physics-based
446 character control. *arXiv preprint arXiv:2308.07200*, 2023.

447 Checklist

- 448 1. For all authors...
- 449 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
450 contributions and scope? [Yes] We provide the environments, quantitative and
451 qualitative results on them in our main paper and supplement.
- 452 (b) Did you describe the limitations of your work? [Yes] Yes, in Sec. 6
- 453 (c) Did you discuss any potential negative societal impacts of your work? [Yes] Yes, in
454 supplement.
- 455 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
456 them? [Yes] Yes.
- 457 2. If you are including theoretical results...
- 458 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 459 (b) Did you include complete proofs of all theoretical results? [N/A]
- 460 3. If you ran experiments (e.g. for benchmarks)...
- 461 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
462 imental results (either in the supplemental material or as a URL)? [Yes] Code and
463 environment will be included in the supplement and open-sourced.
- 464 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
465 were chosen)? [Yes] Yes, in the supplement.
- 466 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
467 ments multiple times)? [Yes] We report our result averaging 1024 env runs.
- 468 (d) Did you include the total amount of computing and the type of resources used (e.g.,
469 type of GPUs, internal cluster, or cloud provider)? [Yes] Yes, in sec. 5
- 470 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 471 (a) If your work uses existing assets, did you cite the creators? [Yes] In supplement.
- 472 (b) Did you mention the license of the assets? [Yes] In supplement.
- 473 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
474 In supplement.
- 475 (d) Did you discuss whether and how consent was obtained from people whose data you're
476 using/curating? [N/A] We do not release a dataset.

- 477 (e) Did you discuss whether the data you are using/curating contains personally identifiable
478 information or offensive content? [N/A] We do not release a dataset.
- 479 5. If you used crowdsourcing or conducted research with human subjects...
- 480 (a) Did you include the full text of instructions given to participants and screenshots, if
481 applicable? [N/A] We do not involve participants.
- 482 (b) Did you describe any potential participant risks, with links to Institutional Review
483 Board (IRB) approvals, if applicable? [N/A] We do not involve participants.
- 484 (c) Did you include the estimated hourly wage paid to participants and the total amount
485 spent on participant compensation? [N/A] We do not involve participants.