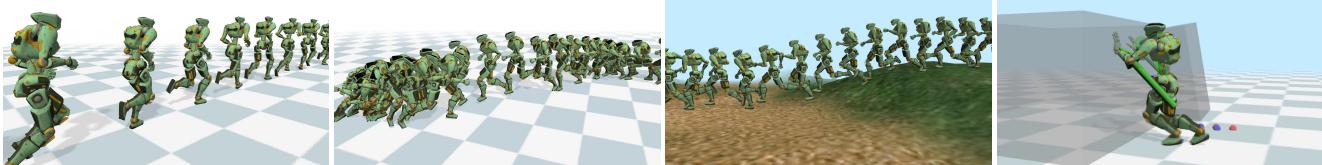


1 Adaptive Tracking of a Single-Rigid-Body Character in Various 2 Environments

3
4 ANONYMOUS AUTHOR(S)
5 SUBMISSION ID: PAPERS_523S1



6
7 Fig. 1. Our character can perform tasks in various environments using a policy learned from a short reference motion clip on flat ground. The results suggest
8 that our controller can perform a range of activities, such as transitioning, climbing on rough terrain, and pushing objects, without additional training.

9 Since the introduction of DeepMimic [Peng et al. 2018], subsequent research
10 has focused on expanding the repertoire of simulated motions across various
11 scenarios. In this study, we propose an alternative approach for this goal,
12 a deep reinforcement learning method based on the simulation of a single-
13 rigid-body character. Using the centroidal dynamics model (CDM) to express
14 the full-body character as a single rigid body (SRB) and training a policy to
15 track a reference motion, we can obtain a policy that is capable of adapting
16 to various unobserved environmental changes and controller transitions
17 without requiring any additional learning. Due to the reduced dimension
18 of state and action space, the learning process is sample-efficient. The final
19 full-body motion is kinematically generated in a physically plausible way,
20 based on the state of the simulated SRB character. The SRB simulation is
21 formulated as a quadratic programming (QP) problem, and the policy outputs
22 an action that allows the SRB character to follow the reference motion. We
23 demonstrate that our policy, efficiently trained within 30 minutes on an
24 ultraportable laptop, has the ability to cope with environments that have
25 not been experienced during learning, such as running on uneven terrain
26 or pushing a box, and transitions between learned policies, without any
27 additional learning.

28 CCS Concepts: • Computing methodologies → Physical simulation.

29 **ACM Reference Format:**

30 Anonymous Author(s). 2023. Adaptive Tracking of a Single-Rigid-Body Char-
31 acter in Various Environments. *ACM Trans. Graph.* 1, 1 (May 2023), 10 pages.
32 <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

33 1 INTRODUCTION

34 Approaches based on deep reinforcement learning (DRL) are pro-
35 ducing significant results in the control of physically simulated
36 characters in recent years. DeepMimic, regarded as one of the pio-
37 neering studies, demonstrated exceptional motion quality for com-
38 plex motions through motion capture data imitation [Peng et al.

39 Permission to make digital or hard copies of all or part of this work for personal or
40 classroom use is granted without fee provided that copies are not made or distributed
41 for profit or commercial advantage and that copies bear this notice and the full citation
42 on the first page. Copyrights for components of this work owned by others than ACM
43 must be honored. Abstracting with credit is permitted. To copy otherwise, or republish,
44 to post on servers or to redistribute to lists, requires prior specific permission and/or a
45 fee. Request permissions from permissions@acm.org.

46 © 2023 Association for Computing Machinery.
47 0730-0301/2023/5-ART \$15.00
48 <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

49 2018]. Expanding the repertoire of motions in diverse scenarios
50 beyond mere imitation has been a focus of subsequent research.
51 Numerous efforts have been made in this direction, such as em-
52 ploying large motion datasets [Bergamin et al. 2019] and generative
53 models [Yao et al. 2022], or training controllers to fulfill a broad
54 range of requirements without relying on specific reference mo-
55 tions [Lee et al. 2021]. However, these approaches often required
56 longer learning time due to the need to experience and adapt to
57 various changes in the environment or tasks.

58 In this study, we propose an alternative and orthogonal approach
59 to the goal of going beyond examples. We express the character as a
60 single rigid body (SRB) by using a simplified physical model called
61 the centroidal dynamics, and train a policy for the SRB character
62 to track a reference motion. Then the learned policy is capable of
63 adapting to various unobserved changes in the environment and
64 controller transitions without any additional learning. Additionally,
65 due to the greatly reduced volume of state and action space, the
66 learning process is sample-efficient. The physical simulation and
67 policy learning are conducted on the SRB character, while final
68 full-body motion is kinematically generated from the simulated SRB
69 states. The policy obtains inputs on the state of the SRB character
70 and produces actions to follow the given reference motion. In the
71 simulation stage, a quadratic programming (QP) solver calculates the
72 contact forces that best achieves the desired acceleration computed
73 from the action and the reference motion. This is then used to update
74 the state of the SRB character.

75 Our approach defines reinforcement learning (RL) tasks in a less
76 specific (excluding full-body details) and more general form (based
77 on the first principle that character movement is caused by con-
78 tact forces), thereby allowing the agent to be less dependent on
79 a particular situation of the full-body character. As a result, our
80 method robustly performs transitions by switching or blending of
81 learned policies, and shows the ability to cope with environments
82 that have not been experienced during learning, such as running on
83 uneven terrain, pushing a box, or balancing against external forces,
84 without any additional learning or parameterization. Furthermore,
85 our method is sample-efficient enough to obtain such a adaptive
86 tracking policy in 30 minutes on an ultraportable laptop.

115 2 RELATED WORK

116 Earlier physics-based motion generation research manually created
 117 a locomotion control algorithm based on error feedback [Coros
 118 et al. 2010; Lee et al. 2010; Yin et al. 2007]. Parameter optimization
 119 and simple balancing rules were used to generate a wide range of
 120 robust motion repertoire [Agrawal et al. 2013; Ha and Liu 2014;
 121 Wang et al. 2012]. Controllers based on QP leveraging the equations
 122 of motion were developed to enhance robustness [Abe et al. 2007;
 123 Da Silva et al. 2008; Kwon and Hodges 2017]. Non-linear trajectory
 124 optimization was used to synthesize physically probable motions
 125 for various tasks by simultaneously considering multiple frames,
 126 as opposed to the single frame optimization used in QP [Mordatch
 127 et al. 2012; Wampler et al. 2014; Ye and Liu 2010]. Model-predictive
 128 control (MPC) methodologies performed online optimization, al-
 129 lowing the construction of complex motions in unknown environ-
 130 ments [Hämäläinen et al. 2015; Macchietto et al. 2009; Tassa et al.
 131 2012]. However, these MPC techniques that utilized the full-body
 132 physical model state space had limitations in terms of motion qual-
 133 ity and long-term planning. MPC techniques with simplified models
 134 allow interactive controls in more complicated environments, fo-
 135 cusing on locomotion tasks [Kwon et al. 2020; Winkler et al. 2018].
 136 Some robotics studies have leveraged simplified models to enhance
 137 their locomotion control strategies [Tsounis et al. 2020; Viereck and
 138 Righetti 2021].

139 RL has been widely used to develop controllers for simulated char-
 140 acters. Earlier studies showed the advantage of creating controllers
 141 with minimal manual effort by designing simple rewards [Coros
 142 et al. 2009; Peng et al. 2015]. DRL expanded the capabilities to tackle
 143 various tasks [Brockman et al. 2016; Duan et al. 2016; Liu and Hod-
 144 gins 2017; Peng et al. 2016; Rajeswaran et al. 2017]. Some studies
 145 employed simplified models or multi-level learning for efficient
 146 long-term planning [Peng et al. 2017; Reda et al. 2022]. Notably,
 147 DeepMimic has achieved impressive motion quality for numerous
 148 reference motions [Peng et al. 2018], but faced limitations in con-
 149 troller performance beyond reference motions. Recent research has
 150 employed large motion datasets and encoding techniques to diver-
 151 sify the motions generated by RL-based controllers and to improve
 152 their generalization capability [Bergamin et al. 2019; Chentanez et al.
 153 2018; Park et al. 2019; Peng et al. 2022; Won et al. 2022; Yao et al.
 154 2022]. Some studies trained controllers to satisfy a wide range of
 155 requirements without using corresponding reference motions, for
 156 example, transitions between different actions or changes in jump
 157 height [Lee et al. 2022, 2021; Peng et al. 2021; Xie et al. 2020; Yin et al.
 158 2021]. Since these full-body simulation-based RL approaches require
 159 observation of various changes in environments and intent during
 160 the learning process, they often require a longer time for learning. In
 161 our method, a policy trained to track only a single reference motion
 162 within a relatively short time can respond to various environmental
 163 changes without experiencing such cases at training time, due to
 164 the flexibility of SRB-based modeling. However, not based on full-
 165 body dynamics, our approach does not directly simulate each part
 166 of the body. As a result, it has the drawback of not being able to
 167 represent contact-rich physical interactions occurring at various
 168 parts of the actual full body. Contact forces can be applied only to
 169 the manually-classified contact points of the SRB character.

170 There have been numerous studies that use a large amount of
 171 motion data to train models and generate realistic real-time full-
 172 body motions purely kinematically [Cho et al. 2021; Holden et al.
 173 2017; Ling et al. 2020; Starke et al. 2019; Zhang et al. 2018]. Our
 174 approach can also be viewed as a kinematic controller, as it generates
 175 the final full-body motions kinematically. However, unlike these
 176 methods that fail to exhibit physical interactions in response to
 177 environmental changes, our approach performs accurate physics
 178 simulation at the level of the SRB. This allows us to create physically
 179 plausible full-body motions by reflecting the changes in the state of
 180 the SRB due to variations in the physical environment.

181 Among various previous studies, the following two studies are
 182 the closest to our study. The research by Xie et al. [Xie et al. 2022]
 183 has many commonalities with our work in its use of the SRB model,
 184 QP, and RL. However, unlike [Xie et al. 2022] which adopted to
 185 generate four-legged locomotion by using manually set gait param-
 186 eters such as foot phase offsets, we focus on creating dynamic and
 187 natural motions of a two-legged character, walking, and various
 188 other motions, by tracking motion capture data. Whereas a simple
 189 Raibert-style heuristic was used to determine foot placement in [Xie
 190 et al. 2022], our policies are trained to output desired foot landing
 191 positions to ensure balancing.

192 [Kwon et al. 2020] is closed to our study in terms of using the
 193 SRB model for a two-legged character, but has the following key
 194 differences: i) [Kwon et al. 2020] generates motions through per-
 195 segment trajectory optimization, which takes too much time for
 196 smooth real-time performance. Our RL-based system can generate
 197 motion at a much faster speed than real-time. ii) To mitigate the
 198 runtime performance issue, [Kwon et al. 2020] proposed a super-
 199 vised learning network that takes pendulum and footstep plans and
 200 generates full-body motion. However, it may fail to produce ade-
 201 quate motion in scenarios deviating significantly from the training
 202 data, such as unexpected external forces. Our method can train an
 203 adaptive controller without additional learning, even in significantly
 204 different scenarios from the reference motion used in training. iii)
 205 In [Kwon et al. 2020], motion generation occurs per-segment, with
 206 planners generating trajectories at half-cycle intervals for a short
 207 future horizon. Consequently, if an unexpected external force is ap-
 208 plied during runtime, the character may respond with a half-cycle
 209 delay. Our policy produces the desired contact position every frame,
 210 enabling an immediate response to external forces.

211 3 SRB CHARACTER AND FRAMES

212 **SRB character** is a simplified representation of
 213 the key physical characteristics of an articulated
 214 character, consisting of a box-shaped rigid body
 215 (the gray box in Figure 3) that approximates the
 216 inertia of the full-body character, as well as four
 217 contact points (the blue spheres) attached on ei-
 218 ther side of the two feet (the black spheres indi-
 219 cating their centers). The orientation of the box
 220 is defined as the orientation of its center of mass,
 221 which implies the overall orientation of the char-
 222 acter. The mass and inertia of the box are set as
 223 the mass (60 kg) of the reference character and its
 224



225 Fig. 3. SRB
 226 character

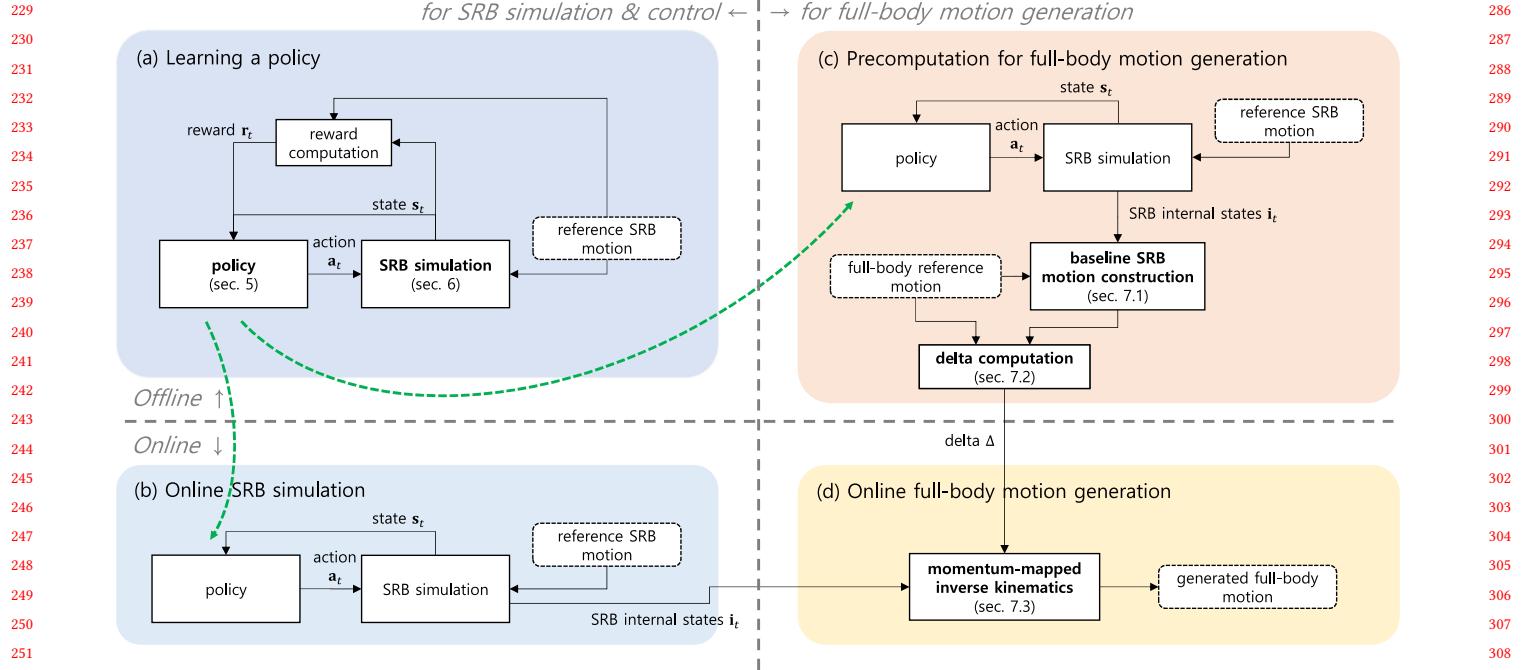


Fig. 2. System overview.

composite rigid body inertia calculated from the default posture (at attention posture).

Reference SRB motion is a reference motion expressed by the SRB character, whose center of mass position and orientation are set to those of the pelvis in the full-body reference motion. The contact timing and position of the feet of the reference SRB motion are set identically to those of the full-body reference motion. The simulated SRB character is controlled to contact the ground at the same time as the reference contact timing.

SRB frames. This paper uses the following three types of reference frames to express various states values:

- **SRB frame** has the SRB character's center of mass as its origin and is attached to the SRB character.
- **Forward-facing SRB frame** has the same origin as the SRB character frame, with its z-axis representing the horizontally projected z-axis of the SRB character frame (forward-facing direction), and its y-axis aligns with the global vertical axis.
- **Projected SRB frame** is obtained by vertically projecting the origin of the forward-facing SRB frame onto the ground or terrain.

Frames corresponding to the reference SRB motion shall be referred to by prefixing "reference" (e.g., "reference SRB frame").

4 OVERVIEW

Figure 2 outlines the overview of the system presented in this paper.

First, a policy is learned for the simulated SRB character to track a reference SRB motion (Figure 2 (a)). In the SRB simulation, the desired acceleration of the SRB character is determined by both

the action from the policy and the reference SRB motion. Then, a quadratic programming (QP) problem is solved to find the actual implementable acceleration that can closely match the desired acceleration based on the given internal state of the SRB character. The internal state of the SRB character is updated by integrating the obtained acceleration. The reward is calculated by comparing the updated SRB internal state and the reference SRB motion.

The learned policy can be used to simulate the SRB character at runtime (Figure 2 (b)). It allows not only tracking a single reference motion but also transitions between different motions through policy switching or blending, without requiring additional learning.

To precompute the data required for full-body motion reconstruction (Figure 2 (c)), the baseline SRB motion, which is the tracking result of the reference SRB motion by the learned policy, is first constructed. Then the difference (Δ) between baseline SRB motion and the full-body reference motion is calculated.

The simulated SRB motion generated at runtime is converted into full-body motion by momentum-mapped inverse kinematics (Figure 2 (d)). The precomputed Δ is combined to produce a more realistic full-body motion by adding the detailed full-body data that were not expressed by the SRB motion.

5 POLICY REPRESENTATION

State $s \in \mathbb{R}^{21}$ is composed of the following components.

- **Center of mass state** $s_c \in \mathbb{R}^{11}$ consists of the height, orientation (in unit quaternion), and generalized velocity (angular and linear velocities) of the center of mass of the SRB character, which are all expressed in the projected SRB frame.

- **Foot state** $s_f^j \in \mathbb{R}^4 (j \in \{\text{left, right}\})$ consists of the center positions and rotations of each foot of the SRB character in the horizontal plane. These values are expressed in the forward-facing SRB frame.
- **Motion phase** $\psi(t) \in [0, 2\pi]$ indicates what point of the reference motion the current motion of the SRB character corresponds. This is actually stored in state s as a 2D-vector $(\sin(\psi(t)), \cos(\psi(t)))$.

Action $a \in \mathbb{R}^{10}$ consists of the following components.

- **Desired foot landing position** $a_s^j \in \mathbb{R}^2 (j \in \{\text{left, right}\})$ in the forward-facing SRB frame. Note that the vertical position is not included.
- **Desired velocity of center of mass** $a_v \in \mathbb{R}^6$ refers to the desired relative linear and angular velocities of the center of mass of the SRB character against that of the reference SRB motion. This is expressed with respect to the SRB frame.

Reward r_t at each step t is configured as below:

$$r_t = w^s r_t^s - w^m r_t^m. \quad (1)$$

Here, r_t^s is the alive reward given when the episode does not end and r_t^m is the mimic reward term that guides the SRB character to closely follow the reference SRB motion.

The mimic reward term r_t^m is calculated as follows:

$$r_t^m = w^P r_t^P + w^e r_t^e. \quad (2)$$

r_t^P is a posture reward term, which gives higher rewards when the states of the centers of mass of both SRB character and reference SRB motion are alike, which is calculated as follows:

$$\begin{aligned} r_t^P = & w^{P_1} \|\Delta p_t - \Delta \hat{p}_{\psi(t)}\| + w^{P_2} \|\Delta R_t \ominus \Delta \hat{R}_{\psi(t)}\| \\ & + w^{P_3} \|p_t - \hat{p}_{\psi(t)}\| + w^{P_4} \|R_t \ominus \hat{R}_{\psi(t)}\|. \end{aligned} \quad (3)$$

Here, Δp_t and ΔR_t respectively refer to the changes in the center of mass' position and orientation of the SRB character between time t and $t+1$, while $\Delta \hat{p}_{\psi(t)}$ and $\Delta \hat{R}_{\psi(t)}$ refer to those changes in the reference SRB motion between the corresponding time points. p_t and R_t respectively refer to the position and orientation of the SRB character's center of mass at time t , while $\hat{p}_{\psi(t)}$ and $\hat{R}_{\psi(t)}$ refer to those of the reference SRB motion. $\|R_1 \ominus R_2\|$ refers to the minimum angle of rotation between the two rotation matrices. To compare the difference in the character height and the leaning angle of the body, the above terms are expressed in the projected frame of the respective character, that is, the projected SRB frame or the reference projected SRB frame.

r_t^e is an end-effector reward, which has a higher value when the support foot positions of both the SRB character and reference SRB motion are alike. This reward is only applied to the foot in contact with the ground in the reference SRB motion and is calculated as follows:

$$r_t^e = \sum_j \sum_k \|c_t^{jk} - \hat{c}_{\psi(t)}^{jk}\|^2, \quad (4)$$

where c_t^{jk} refers to the position of the k^{th} contact point ($k \in \{\text{toe, heel}\}$) of the j^{th} foot ($j \in \{\text{left, right}\}$) of the SRB character at time t with respect to the SRB frame, while $\hat{c}_{\psi(t)}^{jk}$ refers to that of the corresponding contact point of the reference SRB motion at

the corresponding time, expressed with respect to the reference SRB frame. Note that the swing foot is excluded from this term because the SRB character does not directly depict the full-body character's actual swing motion of the foot.

6 SRB SIMULATION

This section will explain the internal state of the SRB character at time t , where the input s_t for the RL policy is extracted, followed by the simulation process by QP that computes the internal state at time $t+1$.

6.1 Internal state of SRB character

The internal state of the SRB character consists of the center of mass' position and orientation, their time derivatives, and the position and orientation of the two feet.

Center of mass state. The global position $p_t^{\{g\}}$ and global orientation $R_t^{\{g\}}$ of the SRB character's center of mass is used to calculate the SRB frame $T_t \in \text{SE}(3)$ at the center of mass as follows:

$$T_t = F(R_t^{\{g\}}, p_t^{\{g\}}), \quad (5)$$

where F is a function that outputs the rigid-body transformation from the position and rotation information.

The time derivatives of the center of mass position and orientation are expressed as the following generalized velocity $\dot{q}_t \in \text{se}(3)$, which indicates the spatial velocity in body frame:

$$[\dot{q}_t] = T_t^{-1} \cdot \dot{T}_t, \quad (6)$$

where \dot{T}_t is the time derivative of T_t , and $[\cdot]$ is the operator that converts $\text{se}(3)$ in the 6-vector coordinates into a 4x4 matrix format.

Foot state consists of the global position f_t^j and orientation F_t^j of each foot. Their specific meaning varies depending on whether the j^{th} foot is in the swing state, determined by the touch down and off timings of that foot in the reference SRB motion. f_t^j is a horizontal position and F_t^j is expressed as only the rotational component against the vertical axis.

When the j^{th} foot is in a swing state, f_t^j moves continuously to the desired foot landing position a_s^j . This a_s^j is not directly suitable as a basis for reproducing the continuous movement of the swing foot since it changes discontinuously over time. Therefore, we use a LQR filter presented in [Hwang et al. 2017], to compute continuous f_t^j from the discontinuous a_s^j . Similarly, by the LQR filter, F_t^j moves continuously to the desired foot landing orientation which is directly obtained from the reference SRB motion. Further details of LQR filtering are described in the supplementary material.

When the j^{th} foot is in a contact state, f_t^j refers to the actual global position of the foot, which remains fixed for the contact duration. The moment the j^{th} foot changes its state from swing to contact, the position of the foot and its rotation against the vertical axis (at the very last moment in the swing phase) become f_t^j and F_t^j for the contact duration. Action a_s^j , which is output from the policy for the foot in the contact state, is ignored and not used.

457 6.2 QP simulation

458 The QP simulator takes the SRB character's center of mass state T_t ,
 459 \dot{q}_t , its relative desired velocity a_v , and contact point c_t as inputs
 460 to calculate its internal state in the next time step. To this end, the
 461 desired acceleration is calculated first based on the reference SRB
 462 motion and a_v . The actual achievable acceleration that maximally
 463 satisfies this desired acceleration and the corresponding contact
 464 force are calculated from the QP formulated as:

$$465 \min_{\dot{q}, \lambda} Q(\dot{q}, \lambda) \quad (7)$$

466 such that

$$467 M\ddot{q} + b = J_f^T F_e + J_c^T B \lambda, \quad (8)$$

$$468 \lambda \geq 0. \quad (9)$$

469 Here, \ddot{q} refers to the generalized center of mass acceleration, λ
 470 : the coefficient for the friction cone basis, F_e : the external force
 471 applied on the SRB character, J_c : the Jacobian matrices related to the
 472 velocity of the contact point c_t , J_f : Jacobian matrices related to the
 473 velocity of the point where the external force is applied, and B : the
 474 friction cone basis vectors. F_e is set to 0 within the RL process, but is
 475 given non-zero values when external forces are applied during the
 476 runtime simulations. In short, the acceleration and contact force that
 477 minimizes the objective function Q are obtained while satisfying the
 478 equation of motion constraints and the constraints on the linearized
 479 basis of contact force [Ellis et al. 2007; Kwon and Hodgins 2017].

480 The objective function Q is defined as follows:

$$481 Q = \|\ddot{q} - \ddot{q}_d\|^2 + w_\lambda \|\lambda\|^2. \quad (10)$$

482 Here, \ddot{q}_d refers to the desired acceleration. w_λ refers to the weight
 483 for the contact force term, which was set to 0.001 to ensure robust
 484 control during our experiments.

485 Action a_v is used in the following to calculate the desired velocity
 486 \dot{q}_d , through which the desired acceleration \ddot{q}_d is obtained.

$$487 \dot{q}_d = \hat{q}_{\psi(t)} + a_v, \quad (11)$$

488 where $\hat{q}_{\psi(t)}$ indicates the linear and angular velocities of the refer-
 489 ence SRB motion's center of mass, and action a_v refers to the desired
 490 relative linear and angular velocities of the SRB character's center
 491 of mass against the reference SRB motion's center of mass.

492 The desired acceleration \ddot{q}_d is calculated as follows by taking the
 493 difference between the desired and current position and velocity.

$$494 \ddot{q}_d = a \log \left(T_t^{-1} \hat{T}_{\psi(t)} \right) + b (\dot{q}_d - \dot{q}_t), \quad (12)$$

495 where T_t is the current SRB frame, while $\hat{T}_{\psi(t)}$ is the reference
 496 SRB frame at the corresponding phase. \dot{q}_t is the current general-
 497 ized velocity, while \dot{q}_d is the desired generalized velocity. The log
 498 function converts a rigid body transformation matrix $\in \text{SE}(3)$ into
 499 a generalized velocity $\in \text{se}(3)$. In our experiments, PD gains a and
 500 b were set at 120 and 35, respectively.

501 By integrating the calculated \ddot{q} , T_t and \dot{q}_t are updated, and the
 502 foot state f_t^j and F_t^j are then updated as described in Section 6.1.
 503 The next time step's state s_{t+1} can be derived from these values.
 504 The simulation process through the QP solver helps the SRB charac-
 505 ter effectively find the contact force that best achieves the desired
 506 acceleration at each moment, helping to learn a robust policy.

514 6.3 Motion phase adjustment

515 At each SRB simulation timestep, the rate at which motion phase
 516 $\psi(t)$ changes is adjusted in two aspects. First, the phase change
 517 rate is increased with increasing locomotion speed to prevent exces-
 518 sively long strides, which is inspired by [Kwon and Hodgins 2017].
 519 Second, when the character experiences a large unexpected force,
 520 it may deviate from the specified contact timings in the reference
 521 SRB motion, resulting in unstable control. To address this, the phase
 522 rate is decreased to delay touchdown when it is expected to hap-
 523 pen before the specified time, and increased to allow for an earlier
 524 touchdown when it is expected to happen later. Further details of
 525 these adjustments are provided in the supplementary material.

526 7 FULL-BODY MOTION GENERATION

527 To generate a full-body motion from a simulated SRB character
 528 motion, we create a baseline SRB motion and calculate the kinematic
 529 and dynamic differences compared to the full-body reference motion.
 530 During real-time simulation, these precomputed differences are
 531 applied to the simulated SRB motion, obtaining target values for
 532 momentum-mapped inverse kinematics (MMIK) to generate the
 533 full-body motion by solving MMIK.

534 7.1 Baseline SRB motion construction

535 After obtaining a trained policy, it is used to simulate the SRB char-
 536 acter to generate multiple simulated cycles, serving as the baseline
 537 SRB motion. Then the COM trajectory and the feet states of the
 538 baseline SRB motion are transformed offline. This transforma-
 539 tion aligns the horizontal start and end points of the COM trajectory with
 540 those of the reference COM trajectory per motion cycle (Figure 6),
 541 while ensuring no inclination of the trajectory by considering only
 542 the rotation against the y-axis.

543 7.2 Delta computation

544 The delta between the baseline SRB motion and the full-body ref-
 545 erence motion is computed to incorporate the reference full-body
 546 details into a reconstructed full-body motion during runtime. The
 547 average differences in the center of mass frame, foot contact points,
 548 centroidal velocity between them are calculated for each different
 549 motion phase and stored as delta Δ .

550 7.3 Momentum-mapped inverse kinematics

551 During runtime, the simulated SRB motion is transformed into
 552 the full-body motion using momentum-mapped inverse kinematics
 553 (MMIK) presented in [Kwon et al. 2020]. We improve the previous
 554 form of MMIK to use the additional input of the precomputed delta
 555 Δ . The inputs are the full-body reference motion, Δ , and simulated
 556 SRB states and it outputs the full-body pose closest to the given full-
 557 body reference pose, while satisfying the footstep position, COM
 558 configuration, and its time derivative of the error-adjusted simulated
 559 SRB motion using Δ . Further details for computing Δ and MMIK
 560 are described in the supplementary material.

561 8 EXPERIMENTAL RESULTS

562 All experiments and policy training were conducted on a ultra-
 563 portable laptop that runs on 16GB RAM, M1 CPU. All motions were

simulated at 60 Hz. Most controllers converged within 30 minutes, around 3M time steps, which is equivalent to about 14 hours of simulated time. The most challenging *Sprint jumps* controller converged within 1-hour, and even after just 30 minutes, a policy performing the motion reliably was obtained. Our approach enables the character to perform a range of reference motions, including those involving drastic changes in direction, such as *Sprint* or *Sharp turns*, or those requiring strong force (Figure 7).

We demonstrate how effectively our controllers respond to diverse, unobserved changes. Note that all controllers used in the experiments were trained simply to mimic reference motions on flat ground (except Section 8.4), without any additional training for adapting to those changes. In some experiments, we conducted comparisons with DeepMimic [Peng et al. 2018]. DeepMimic has the advantage of simulating full-body dynamics, but it is designed to imitate detailed poses of a full-body reference motion, limiting its adaptability to unfamiliar situations. As a result, DeepMimic exhibited less adaptability compared to our method. However, it is important to note that DeepMimic and our method belong to different categories and have different purposes. These comparative experiments were not intended to claim that our method is always superior to DeepMimic, but rather to showcase the advantages of our simplified physics-based approach compared to widely used full-body-based methods. Further implementation and training details, and additional experimental results such as sample efficiency are provided in the supplementary material.

8.1 Controller transitions and interpolations

Transitions by switching. Our method allows for transitioning between different controllers simply by switching pre-trained policies (Figure 8). At the transition point, the policy π_A , reference SRB motion \hat{m}_A , and precomputed delta Δ_A of the Controller A are immediately replaced with those of the Controller B. We use a motion stitching technique that gradually reduces the difference between the reference motions and deltas at the transition point and reflects it over time. Note that the transition was successful (no falling within 20 seconds after the switch) when the difference in movement speed was within 30%, such as *Sprint* and *Sprint jumps* or *Sprint* and *Sharp turns*, or when both controllers are exceptionally robust (e.g. *Sprint* and *Run*).

Transitions by blending. Transitions between controllers with larger speed differences can be achieved through blending, where Controller A gradually transitions to Controller B by linearly increasing the interpolation weight t from 0 to 1 (for 1 sec). The actions from π_A and π_B , \hat{m}_A and \hat{m}_B , Δ_A and Δ_B , as well as the change rates of $\psi_A(t)$ and $\psi_B(t)$, the remaining time in \hat{m}_A and \hat{m}_B until each j^{th} foot touches down $t_A^{d,j}, t_B^{d,j}$ and until touch off $t_A^{o,j}, t_B^{o,j}$ are interpolated. With this method, smooth transitions were achieved between controllers with significant differences in movement speeds and styles, such as *Fast walk* (1.8 m/s) and *Run* (3.6 m/s) (Figure 9).

Interpolation. The interpolation method allows for creating new interpolated controllers at constant ratios between two controllers, apart from transitioning between two controllers (Figure 10).

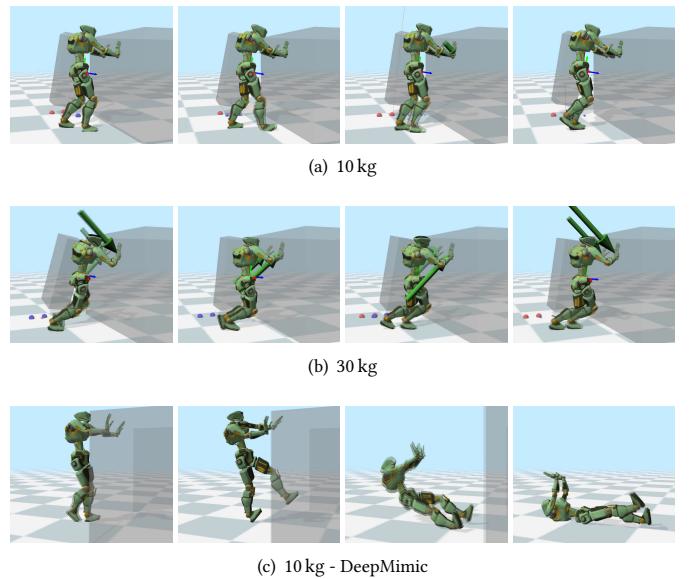


Fig. 4. Pushing a box of varying weights.

8.2 Environmental adaptations

Uneven terrain. Our controllers trained on even terrain was able to create locomotion on uneven terrain with a gradient of a maximum 30 degrees (Figure 11). This is achieved without the need for any additional modifications to the algorithm, and the results are based on the inherent adaptability of the SRB-based policy. In contrast, DeepMimic was observed to struggle even on gentle slopes, as it focuses on tracking full-body poses, resulting in the swing foot tripping.

Push boxes. The policy was trained to track an edited version of the *Walk* reference motion, where both arms were raised forward. The learning environment did not include any boxes, and therefore, the policy did not receive any box information. As shown in Figure 4, the character was able to push boxes of different weights using different motions, by leaning forward to exert more force when the box is heavier. Adjusting the foot landing position slightly behind the desired landing position in the action output of the policy resulted in a stronger pushing behavior when dealing with heavy boxes. In contrast, a DeepMimic controller trained to track the same reference motion was unable to push the box. Our SRB-based model achieved comparable results to a previous study [Lee et al. 2021] without the need for a policy parameterized by box weights.

8.3 Comparision for external pushes

We compared the balance maintenance level of characters controlled by our method, FFMLCD [Kwon et al. 2020], and DeepMimic [Peng et al. 2018]. As reported in a previous study [Lee et al. 2015], the response to external force can vary depending on the push timing, so we performed push experiments throughout the entire motion phase rather than at a specific point. The experiments used *Sprint* controllers and were performed at 20 evenly spaced phase points.

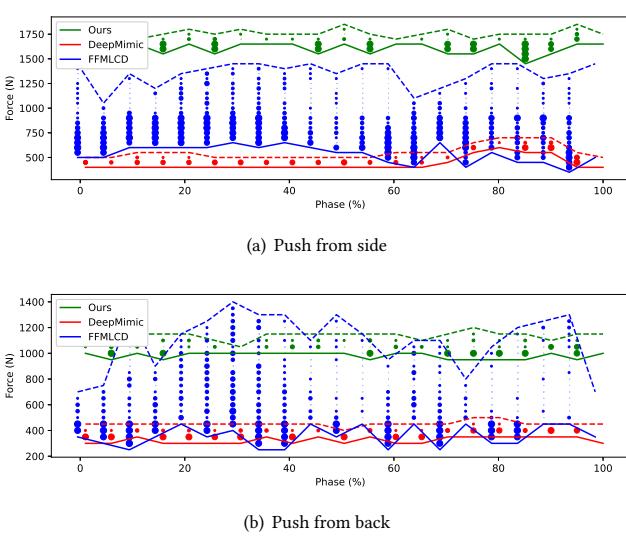


Fig. 5. Plots for external pushes. Dotted line: the smallest external force at which the character always loses balance. Solid line: the largest external force at which the character maintains balance at each phase without falling. Circles: the ratio of successful balance maintenance (no falling within 20 seconds after the push) out of 10 trials at points where both balance maintenance and loss occur, with its size representing the proportion of successful balance maintenance.

External force was applied to the character’s root from the side or behind for 0.2 seconds, ranging from 100 N to 2000 N in 50 N increments. Each force level was applied 10 times at each point.

As depicted in Figure 5, our method shows a higher minimum force at which the character always loses balance (dotted lines) in side pushes, while it is comparable to FFMLCD in back pushes. However, from a practical standpoint, the important factor is the force magnitude at which the character maintains balance without falling (solid lines). In this aspect, our method exhibits a force magnitude approximately three times higher than FFMLCD in both push cases. This is because our method, utilizing RL, generates desired footstep positions at each frame to maintain balance, while FFMLCD relies on non-linear optimization-based planning that may occasionally struggle to find optimal solutions in challenging situations. Moreover, FFMLCD performs this planning at each half cycle instead of each frame, which further limits its responsiveness. Due to the limited adaptability caused by full-body tracking, DeepMimic easily fell even with weaker forces. All controllers were more stable when pushed from the side compared to the back, because pushing from the back would cause the body to lean forward and the swing foot to suddenly hit the ground, resulting in a fall or excessive contact force.

8.4 Interactive control

We also conducted experiments to enable more useful applications by allowing users to interactively control the facing direction of the character. For this purpose, we added a scalar δ_y , representing the difference between the desired facing direction and the current SRB

facing direction, to the inputs of the policy network. The policy is then trained using the following reward r_t^c :

$$r_t^c = \exp(-r_t) \cdot \exp(-w^c \cdot (\cos(\delta_y) - 1)), \quad (13)$$

where r_t represents the reward defined in Equation 1, and w^c denotes the weight.

9 DISCUSSION

In this paper, we have presented a framework that can learn an adaptive tracking policy in a sample-efficient manner using a SRB character simplifying a full-body character. The QP-based SRB simulation allows efficient learning of a policy, and the simulated SRB motion can be converted back into a full-body motion using precomputed delta and a momentum-mapped inverse kinematics solver.

Thanks to the less detailed SRB model-based learning, our method enables fast policy learning with fewer samples compared to traditional full-body methods. Moreover, the learned policy exhibits adaptability in various environments. However, it is worth noting that our method has limitations in accurately representing contact-rich scenarios, unlike full-body physics simulations. For example, it does not explicitly represent collisions between body parts and the ground when the character falls. In our method, "falling" refers to when the center of mass is too close to the ground, making the solution of MMIK meaningless. In practice, we can partially address this limitation by switching to a full-body simulation using PD-servo to generate a more natural fall when the character begins to lose balance.

In our method, each controller is trained based on a single short reference motion clip. Instead of relying on a large number of reference motions, our approach prioritizes maximizing the number of situations that the same controller can handle without additional learning. Once such an adaptability is achieved, we believe that it can be straightforwardly extended in the direction of enabling a wider range of motions by using more complex network structures supporting unorganized motion capture datasets, such as VAE or RNN [Park et al. 2019; Yao et al. 2022]. The proposed framework allows only indirect adjustment of the contact timing based on the adjustment of the speed of phase progression. This can negatively affect robustness because a limited range of contact timing adjustment interferes with more flexible adaptation of the character to external force. However, unlike DeepMimic-style controllers which use fixed contact timing, the contact timing is easily adjustable in our framework that employs virtual legs, and thus it would help to improve the robustness of the controllers by adjusting the contact timing in a similar way as tested in [Reda et al. 2022]. One of the computational bottlenecks in our learning process is solving the QP. As future work, it would be beneficial to explore the use of a GPU-based QP solver, such as the one presented in [Xie et al. 2022], to further accelerate learning.

REFERENCES

- Yeuhi Abe, Marco Da Silva, and Jovan Popović. 2007. Multiobjective control with frictional contacts. In *Proceedings of the 2007 ACM SIGGRAPH/Eurographics symposium on Computer animation*. 249–258.
- Shailesh Agrawal, Shuo Shen, and Michiel van de Panne. 2013. Diverse motion variations for physics-based character animation. In *Proceedings of the 12th ACM SIGGRAPH/Eurographics Symposium on Computer Animation*. 37–44.

- 99 Kevin Bergamin, Simon Clavet, Daniel Holden, and James Richard Forbes. 2019. DReCon:
 800 data-driven responsive control of physics-based characters. *ACM Transactions on Graphics (TOG)* 38, 6 (Nov. 2019), 206:1–206:11. 856
 801 Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman,
 802 Jie Tang, and Wojciech Zaremba. 2016. Openai gym. *arXiv preprint arXiv:1606.01540* 857
 803 (2016). 858
 804 Nuttапong Chentanez, Matthias Müller, Miles Macklin, Viktor Makoviychuk, and Stefan
 805 Jeschke. 2018. Physics-Based Motion Capture Imitation with Deep Reinforcement
 806 Learning. In *Proceedings of the 11th ACM SIGGRAPH Conference on Motion, Interaction
 807 and Games (MIG ’18)*. Article 1, 10 pages. 859
 808 Kyungmin Cho, Chaelin Kim, Jungjin Park, Joonkyu Park, and Junyong Noh. 2021.
 809 Motion recommendation for online character control. *ACM Transactions on Graphics* 40, 6 (2021). 860
 810 Stelian Coros, Philippe Beaudoin, and Michiel van de Panne. 2009. Robust task-based
 811 control policies for physics-based characters. *ACM Trans. Graph. (Proc. SIGGRAPH
 812 Asia)* 28, 5 (2009), 1–9. 861
 813 Stelian Coros, Philippe Beaudoin, and Michiel Van de Panne. 2010. Generalized biped
 814 walking control. *ACM Transactions On Graphics (TOG)* 29, 4 (2010), 1–9. 862
 815 Marco Da Silva, Yeuhi Abe, and Jovan Popović. 2008. Simulation of human motion data
 816 using short-horizon model-predictive control. In *Computer Graphics Forum*, Vol. 27. 863
 817 371–380. 864
 818 Yan Duan, Xi Chen, Rein Houthooft, John Schulman, and Pieter Abbeel. 2016. Bench-
 819 marking deep reinforcement learning for continuous control. In *Proceedings of the
 820 33rd International Conference on International Conference on Machine Learning -
 821 Volume 48 (ICML’16)*. 1329–1338. 865
 822 Jane Ellis, Harald Winkler, Jan Corfee-Morlot, and Frédéric Gagnon-Lebrun. 2007. CDM:
 823 Taking stock and looking forward. *Energy policy* 35, 1 (2007), 15–28. 866
 824 Sehoon Ha and C Karen Liu. 2014. Iterative training of dynamic skills inspired by
 825 human coaching techniques. *ACM Transactions on Graphics (TOG)* 34, 1 (2014),
 826 1–11. 867
 827 Perttu Hämäläinen, Joose Rajamäki, and C Karen Liu. 2015. Online control of simulated
 828 humanoids using particle belief propagation. *ACM Transactions on Graphics (TOG)*
 829 34, 4 (2015), 1–13. 868
 830 Daniel Holden, Taku Komura, and Jun Saito. 2017. Phase-Functioned Neural Networks
 831 for Character Control. *ACM Transactions on Graphics* 36, 4, Article 42 (2017),
 832 13 pages. 869
 833 Jaeyoung Hwang, Kwanguk Kim, Il Hong Suh, and Taesoo Kwon. 2017. Performance-
 834 based animation using constraints for virtual object manipulation. *IEEE computer
 835 graphics and applications* 37, 4 (2017), 95–102. 870
 836 Taesoo Kwon and Jessica K Hodgins. 2017. Momentum-mapped inverted pendulum
 837 models for controlling dynamic human motions. *ACM Transactions on Graphics
 838 (TOG)* 36, 1 (2017), 1–14. 871
 839 Taesoo Kwon, Yoonsang Lee, and Michiel Van De Panne. 2020. Fast and flexible
 840 multilegged locomotion using learned centroidal dynamics. *ACM Transactions on
 841 Graphics (TOG)* 39, 4 (2020), 46–1. 872
 842 Seyoung Lee, Jiye Lee, and Jehee Lee. 2022. Learning Virtual Chimeras by Dynamic
 843 Motion Reassembly. *ACM Transactions on Graphics* 41, 6 (2022), 182:1–182:13. 873
 844 Seyoung Lee, Sunmin Lee, Yongwoo Lee, and Jehee Lee. 2021. Learning a family of
 845 motor skills from a single motion clip. *ACM Transactions on Graphics* 40, 4 (July
 846 2021), 93:1–93:13. 874
 847 Yoongsang Lee, Sungeun Kim, and Jehee Lee. 2010. Data-driven biped control. *ACM
 848 Trans. Graph.* 29, 4 (2010), 1–8. 875
 849 Yoongsang Lee, Kyungho Lee, Soon-Sun Kwon, Jiwon Jeong, Carol O’Sullivan, Moon Seok
 850 Park, and Jehee Lee. 2015. Push-recovery Stability of Biped Locomotion. *ACM
 851 Transactions on Graphics (TOG)* 34, 6 (2015), 180:1–180:9. 876
 852 Hung Yu Ling, Fabio Zimmo, George Cheng, and Michiel Van De Panne. 2020. Character
 853 controllers using motion VAEs. *ACM Transactions on Graphics* 39, 4 (2020). 877
 854 Libin Liu and Jessica Hodgins. 2017. Learning to schedule control fragments for physics-
 855 based characters using deep q-learning. *ACM Transactions on Graphics (TOG)* 36, 3
 856 (2017), 1–14. 878
 857 Adriano Macchietto, Victor Zordan, and Christian R. Shelton. 2009. Momentum control
 858 for balance. *ACM Transactions on Graphics* 28, 3 (July 2009), 80:1–80:8. 879
 859 Igor Mordatch, Emanuel Todorov, and Zoran Popović. 2012. Discovery of complex
 860 behaviors through contact-invariant optimization. *ACM Transactions on Graphics
 861 (TOG)* 31, 4 (2012), 1–8. 880
 862 Soohwan Park, Hoseok Ryu, Seyoung Lee, Sunmin Lee, and Jehee Lee. 2019. Learning
 863 predict-and-simulate policies from unorganized human motion data. *ACM
 864 Transactions on Graphics* 38, 6 (2019), 205:1–205:11. 881
 865 Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. 2018. Deepmimic:
 866 Example-guided deep reinforcement learning of physics-based character skills. *ACM
 867 Transactions on Graphics (TOG)* 37, 4 (2018), 1–14. 882
 868 Xue Bin Peng, Glen Berseth, and Michiel Van de Panne. 2015. Dynamic terrain traversal
 869 skills using reinforcement learning. *ACM Transactions on Graphics (TOG)* 34, 4
 870 (2015), 1–11. 883
 871 Xue Bin Peng, Glen Berseth, and Michiel Van de Panne. 2016. Terrain-adaptive locomotion
 872 skills using deep reinforcement learning. *ACM Transactions on Graphics (TOG)*
 873 35, 4 (2016), 1–12. 884
 874 Xue Bin Peng, Glen Berseth, KangKang Yin, and Michiel van de Panne. 2017. DeepLoco:
 875 Dynamic Locomotion Skills Using Hierarchical Deep Reinforcement Learning. *ACM
 876 Transactions on Graphics (Proc. SIGGRAPH 2017)* 36, 4 (2017). 885
 877 Xue Bin Peng, Yunrong Guo, Lina Halper, Sergey Levine, and Sanja Fidler. 2022. ASE:
 878 large-scale reusable adversarial skill embeddings for physically simulated characters.
 879 *ACM Transactions on Graphics* 41, 4 (July 2022), 94:1–94:17. 886
 880 Xue Bin Peng, Ze Ma, Pieter Abbeel, Sergey Levine, and Angjoo Kanazawa. 2021.
 881 AMP: adversarial motion priors for stylized physics-based character control. *ACM
 882 Transactions on Graphics* 40, 4 (2021), 144:1–144:20. 887
 883 Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman,
 884 Emanuel Todorov, and Sergey Levine. 2017. Learning complex dexterous
 885 manipulation with deep reinforcement learning and demonstrations. *arXiv preprint
 886 arXiv:1709.10087* (2017). 887
 887 Daniele Reda, Hung Yu Ling, and Michiel van de Panne. 2022. Learning to Brachiate
 888 via Simplified Model Imitation. In *ACM SIGGRAPH 2022 Conference Proceedings
 889 (SIGGRAPH ’22)*. Article 24, 9 pages. 888
 890 Sebastian Starke, He Zhang, Taku Komura, and Jun Saito. 2019. Neural State Machine
 891 for Character-Scene Interactions. *ACM Transactions on Graphics* 38, 6, Article 209
 892 (2019), 14 pages. 889
 893 Yuval Tassa, Tom Erez, and Emanuel Todorov. 2012. Synthesis and stabilization of com-
 894 plex behaviors through online trajectory optimization. In *2012 IEEE/RSJ International
 895 Conference on Intelligent Robots and Systems*. IEEE, 4906–4913. 890
 896 Vassilios Tsounis, Mitja Alge, Joonho Lee, Farbod Farshidian, and Marco Hutter. 2020.
 897 Deepgait: Planning and control of quadrupedal gait using deep reinforcement
 898 learning. *IEEE Robotics and Automation Letters* 5, 2 (2020), 3699–3706. 899
 899 Julian Viereck and Ludovic Righetti. 2021. Learning a centroidal motion planner for
 900 legged locomotion. In *2021 IEEE International Conference on Robotics and Automation
 901 (ICRA)*. IEEE, 4905–4911. 900
 902 Kevin Wampler, Zoran Popović, and Jovan Popović. 2014. Generalizing locomotion
 903 style to new animals with inverse optimal regression. *ACM Transactions on Graphics
 904 (TOG)* 33, 4 (2014), 1–11. 901
 905 Jack M Wang, Samuel R Hamner, Scott L Delp, and Vladlen Koltun. 2012. Optimizing
 906 locomotion controllers using biologically-based actuators and objectives. *ACM
 907 Transactions on Graphics (TOG)* 31, 4 (2012), 1–11. 902
 908 Alexander W Winkler, C Dario Bellicoso, Marco Hutter, and Jonas Buchli. 2018. Gait
 909 and trajectory optimization for legged systems through phase-based end-effector
 910 parameterization. *IEEE Robotics and Automation Letters* 3, 3 (2018), 1560–1567. 903
 911 Jungdam Won, Deepak Gopinath, and Jessica Hodgins. 2022. Physics-based character
 912 controllers using conditional VAEs. *ACM Transactions on Graphics* 41, 4 (2022),
 913 96:1–96:12. 904
 914 Zhaoming Xie, Xingye Da, Buck Babich, Animesh Garg, and Michiel van de Panne.
 915 2022. Glide: Generalizable quadrupedal locomotion in diverse environments with
 916 a centroidal model. In *Algorithmic Foundations of Robotics XV: Proceedings of the
 917 Fifteenth Workshop on the Algorithmic Foundations of Robotics*. Springer, 523–539. 905
 918 Zhaoming Xie, Hung Yu Ling, Nam Hee Kim, and Michiel van de Panne. 2020. ALL-
 919 STEPS: Curriculum-driven Learning of Stepping Stone Skills. *Computer Graphics
 920 Forum* 39, 8 (2020), 213–224. 906
 921 Heyuan Yao, Zhenhua Song, Baoquan Chen, and Libin Liu. 2022. ControlVAE: Model-
 922 Based Learning of Generative Controllers for Physics-Based Characters. *ACM
 923 Transactions on Graphics* 41, 6 (2022), 183:1–183:16. 907
 924 Yuting Ye and C. Karen Liu. 2010. Optimal feedback control for character animation
 925 using an abstract model. *ACM Trans. Graph.* 29, 4 (2010), 1–9. 908
 926 KangKang Yin, Kevin Loken, and Michiel Van de Panne. 2007. Simbicon: Simple biped
 927 locomotion control. *ACM Transactions on Graphics (TOG)* 26, 3 (2007), 105–es. 909
 928 Zhiqi Yin, Zeshi Yang, Michiel Van De Panne, and Kangkang Yin. 2021. Discovering
 929 diverse athletic jumping strategies. *ACM Transactions on Graphics* 40, 4 (July 2021),
 930 91:1–91:17. 910
 931 He Zhang, Sebastian Starke, Taku Komura, and Jun Saito. 2018. Mode-Adaptive Neural
 932 Networks for Quadruped Motion Control. *ACM Transactions on Graphics* 37, 4
 933 (2018). 911
 934 935

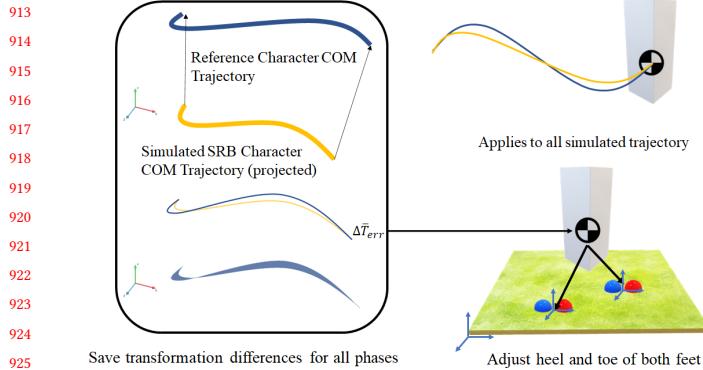


Fig. 6. Conceptual depiction of aligning COM trajectory and calculating COM frame delta.

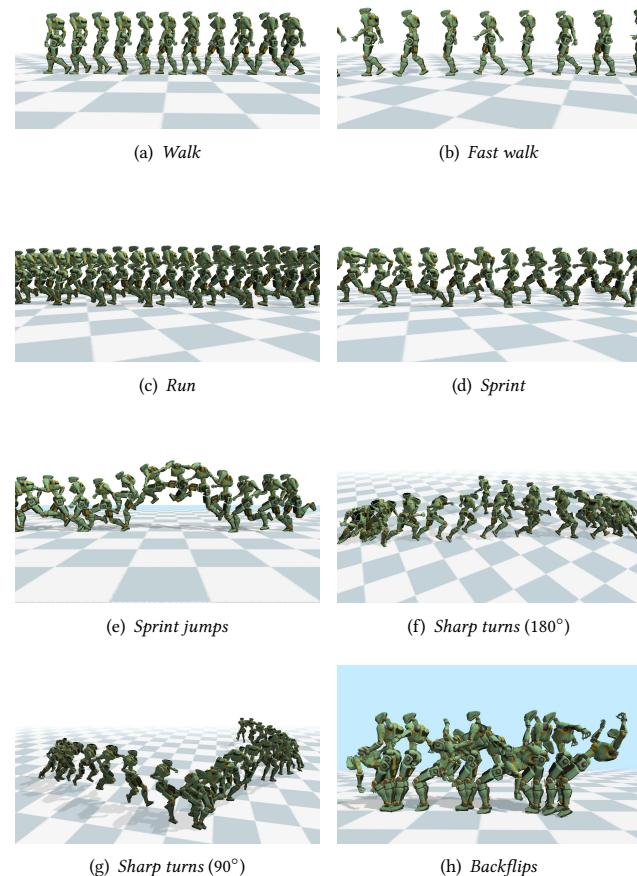


Fig. 7. Tracking various motions.

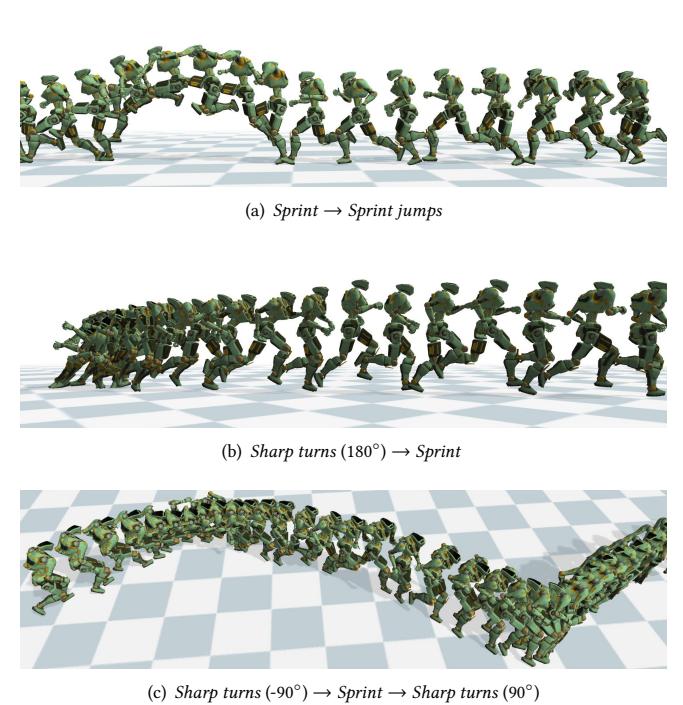


Fig. 8. Transitions by controller switching.



Fig. 9. Blending of controllers. Bi-directional transitions of motions are made between Fast walk (1.8 m/s) and Run (3.6 m/s) through the blended controller.

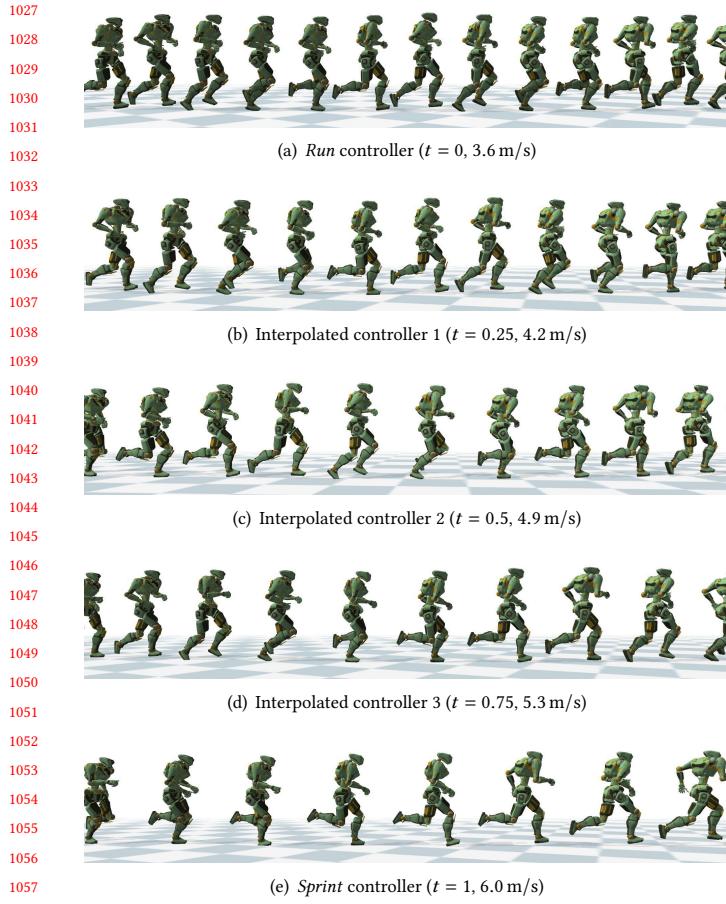


Fig. 10. Interpolated controllers. This figure shows the simulated motions from *Run* and *Sprint* controllers and the interpolated controllers between the two. The speed in parentheses is the moving speed of the simulated character.

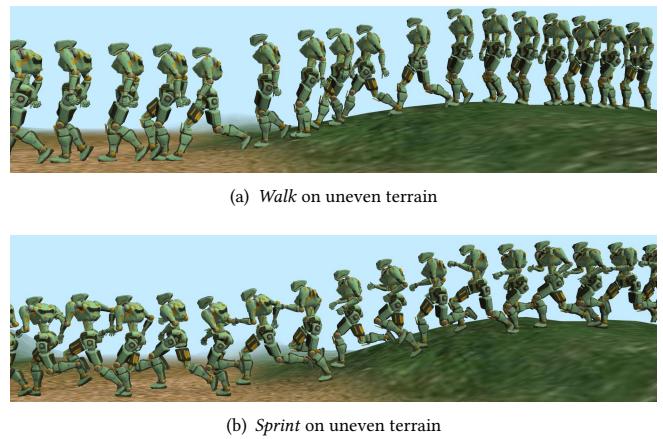


Fig. 11. Adaptation for uneven terrain.