

AdaptNet: Policy Adaptation for Physics-Based Character Control

PEI XU, Clemson University, USA and Roblox, USA

KAIXIANG XIE, McGill University, Canada

SHELDON ANDREWS, École de Technologie Supérieure, Canada and Roblox, USA

PAUL G. KRY, McGill University, Canada

MICHAEL NEFF, University of California, Davis, USA

MORGAN MCGUIRE, Roblox, USA and University of Waterloo, Canada

IOANNIS KARAMOZAS, University of California, Riverside, USA

VICTOR ZORDAN, Roblox, USA and Clemson University, USA

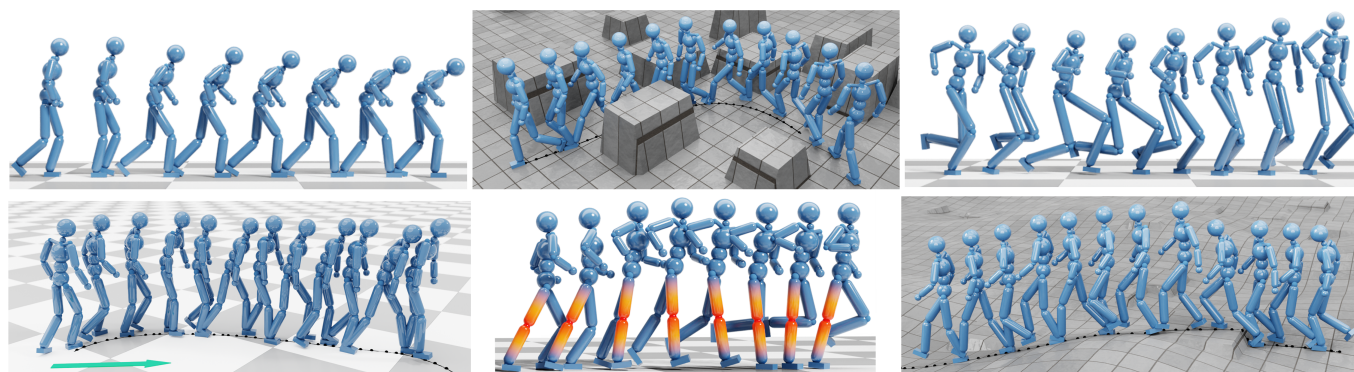


Fig. 1. Examples policy adaptation for locomotion. From left to right and top to bottom: motion interpolation, local collision avoidance, body-length changes, style transfer, morphology changes, rough terrain adaptation.

Motivated by humans' ability to adapt skills in the learning of new ones, this paper presents AdaptNet, an approach for modifying the latent space of existing policies to allow new behaviors to be quickly learned from like tasks in comparison to learning from scratch. Building on top of a given reinforcement learning controller, AdaptNet uses a two-tier hierarchy that augments the original state embedding to support modest changes in a behavior and further modifies the policy network layers to make more substantive changes. The technique is shown to be effective for adapting existing physics-based controllers to a wide range of new styles for locomotion, new task targets, changes in character morphology and extensive changes in environment. Furthermore, it exhibits significant increase in learning efficiency, as indicated by greatly reduced training times when compared to training from

scratch or using other approaches that modify existing policies. Code is available at <https://motion-lab.github.io/AdaptNet>.

CCS Concepts: • **Computing methodologies** → **Animation**; *Physical simulation*; *Reinforcement learning*.

Additional Key Words and Phrases: character animation, physics-based control, motion synthesis, reinforcement learning, motion style transfer, domain adaptation, GAN

ACM Reference Format:

Pei Xu, Kaixiang Xie, Sheldon Andrews, Paul G. Kry, Michael Neff, Morgan McGuire, Ioannis Karamouzas, and Victor Zordan. 2023. AdaptNet: Policy Adaptation for Physics-Based Character Control. *ACM Trans. Graph.* 42, 6, Article 112.1522 (December 2023), 18 pages. <https://doi.org/10.1145/3618375>

1 INTRODUCTION

Research on physically-based character animation has received a great deal of attention recently, especially using reinforcement learning (RL) to develop control policies that produce a wide spectrum of motion behaviors and styles with few or no manual inputs. Most techniques rely on reference human motion to either provide direct tracking or indirect comparison to constrain movement, along with additional targets and rewards to shape task success (e.g., [Liu and Hodgins 2018; Peng et al. 2018a; Xu and Karamouzas 2021]). However, methods to date largely develop policies or controllers for a known behavior, and must be relearned (usually from scratch) to produce a new behavior. While curriculum-style learning and warm-start approaches may be used to migrate policies to targeted

Authors' addresses: Pei Xu, Clemson University, USA and Roblox, USA, peix@clemson.edu; Kaixiang Xie, McGill University, Canada, kaixiang.xie@mail.mcgill.ca; Sheldon Andrews, École de Technologie Supérieure, Canada and Roblox, USA, sheldon.andrews@gmail.com; Paul G. Kry, McGill University, Canada, kry@cs.mcgill.ca; Michael Neff, University of California, Davis, USA, mpneff@ucdavis.edu; Morgan McGuire, Roblox, USA and University of Waterloo, Canada, morgan@roblox.com; Ioannis Karamouzas, University of California, Riverside, USA, ioannis@cs.ucr.edu; Victor Zordan, Roblox, USA and Clemson University, USA, vbzordan@roblox.com.

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

© 2023 Association for Computing Machinery.

0730-0301/2023/12-ART112.1522 \$15.00

<https://doi.org/10.1145/3618375>

goal tasks [Tao et al. 2022; Yin et al. 2021], we instead aim to broadly adapt previously trained policies to make them usable in a wide spectrum of new scenarios without the need for full retraining.

Inspired by recent work in conditioning existing models in image-based stable diffusion and large language models [Hu et al. 2021; Zhang and Agrawala 2023], we introduce *AdaptNet* as an approach for controlling physically based characters that modifies an existing policy to produce behavior in a variety of new settings. The main novelty of our work is the ability to control the motion generation process through editing the latent space. In physics-based character control tasks, there is an opportunity to better understand and exploit the latent space representation of control policies obtained using reinforcement learning frameworks. AdaptNet provides an initial step in this direction.

Specifically, our approach relies on the training of weights for new network components that are injected into a previously trained policy network. Building on top of a pre-existing multi-objective reinforcement learning controller, we propose a two-tier architecture for AdaptNet that augments the latent state embedding while adding modifications to the remaining layers for control refinement. The first layer modifies the latent space projected from the association of the task and character state. It supports adding elements to the control state, as well as changing the imitation and task rewards. Meanwhile, the deeper, control-level refinement augments the policy’s action derived from the latent state, supporting more substantive changes to the task control. Together, AdaptNet performs fast training from a previously trained policy and is capable of making a wide spectrum of adaptations from a single behavior.

As in Figure 1, we showcase our learning framework with numerous controller adaptation examples, including changes in the style of locomotion derived from very short reference motions. AdaptNet can perform this “few-shot style transfer” using only the embedding layer augmentation in a fraction of the time it takes to learn the original locomotion policy. Furthermore, through interpolating in the latent space, it is possible to control the generated control dynamically and smoothly transition from the original behavior to the new style. We further experiment with changes to the character morphology by “locking” joints and changing limb lengths. While these changes lead to failure in the original policy, AdaptNet augments the policy easily to account for the various changes. We also investigate changes in the environment, exploring adaptation for locomotion on rough and slick (low-friction) terrains, as well as on obstacle-filled environments. In each case, AdaptNet provides significant improvement leading to characters that robustly traverse a range of new settings (see Figure 1 and accompanying video).

We evaluate the effectiveness of AdaptNet on various tasks, including its ability for adaptation of imitation learning, different goal rewards, and environmental states. We compare our approach against training from scratch, as well as training-continuation (finetuning). Training with AdaptNet can typically be carried out within 10-30 minutes for simple adaptation tasks, and up to 4 hours for complex locomotion tasks and environment changes. Within such modest training time budgets, in most cases it is impossible to obtain a working controller that can adhere to imitation and goal-task objectives when training from scratch or finetuning a pre-existing policy. Additional ablation studies support the specific architecture

we propose over several alternatives along with highlighting AdaptNet’s ability to successfully control and modify the latent space.

The contributions of our work are summarized as follows:

- We show how the latent space representation of an RL policy can be modified for motion synthesis in physics-based motor control tasks.
- Based on this, we introduce AdaptNet as a framework to efficiently modify a pre-trained physics-based character controller to new tasks.
- We showcase the applicability of AdaptNet on a variety of multi-objective adaptation tasks, including few-shot motion style transfer, motion interpolation, character morphology adaptation, and terrain adaptation.

2 RELATED WORK

Our approach follows a wide set of previous related work stemming from general disciplines in computer animation, robotics, machine learning and image generation. We focus on the background that is most relevant, categorized in physically based character skill control, transfer learning, and latent space adaptation.

2.1 Deep Reinforcement Learning for Skilled Motion

Deep learning neural network control policies have become the staple for physics-based character animation research due to their ability to synthesize a range of skilled motions. In recent years, techniques have trained control policies to animate physics-based humanoid characters for agile motions [Yin et al. 2021], team sports [Liu and Hodgins 2018; Xie et al. 2022], martial arts [Won et al. 2021], juggling [Chemin and Lee 2018; Luo et al. 2021; Xu et al. 2023], performing complex environment interactions [Merel et al. 2020], as well as general locomotion tasks [Bergamin et al. 2019; Peng et al. 2018a]. The recent survey by Kwiatkowski et al. [2022] provides a comprehensive overview of approaches that have been developed for motion synthesis and control of animated characters.

Training skill-specific policies often requires extended training time, necessitating years of simulated learning [Peng et al. 2022]. Skill re-use and combining pre-trained policies to perform more complex tasks offer an alternative that can create needed savings from this extensive training. To this end, a number of papers have proposed ways to reuse and/or combine policies. For example, DeepMimic [Peng et al. 2018a] trains a composite policy that transitions between a collection of different skills. Liu and Hodgins [2017] experiment with hierarchical models that sequence a set of pre-trained control fragments. Hejna et al. [2020] explore a hierarchical approach to decouple low and high-level policies to transfer skills from agents with simple morphologies to more complex ones, and found that it helps to reduce overall sampling. Likewise, we demonstrate that the proposed AdaptNet approach is effective when adapting pre-trained policies to new character morphologies and motion styles with relatively little additional training time.

Curriculum learning is also related to skill adaptation since the agent is trained on tasks with increasing difficulty [Karpthy and van de Panne 2012; Yu et al. 2018]. The approach is demonstrated to be effective for training controllers that allow agents to traverse environments of increasing complexity [Heess et al. 2017; Xie et al.

2020] and recover to standing [Frezzato et al. 2022] under increasingly challenging conditions. In comparison, we demonstrate that our approach efficiently allows a physically simulated humanoid to adapt pre-trained walking and running skills to new terrains as well. However, the aim for curriculum learning is somewhat different than our own in that it is usually used as a means to develop a single advanced skill while we focus on the ability to generalize from one behavior to many.

2.2 Transfer Learning

In machine learning, a common approach for model adaptation is to start with a pre-trained model and fine tune it on a new task. Over the years a number of architectures have been proposed to overcome the overfitting and expressivity issues of finetuning, including GAN-inspired approaches for domain adaptation [Ganin et al. 2016; Tzeng et al. 2017] and adding new models to previously learnt ones through lateral connections [Rusu et al. 2016b, 2017]. To facilitate better model transfer, algorithms have been explored that account for entropy optimization [Haarnoja et al. 2017; Wang et al. 2021]. As well, others directly manipulate the source task domain through randomizing physical parameters of the agent and/or environment while adapting the source domain to the target one [Ganin et al. 2016; Peng et al. 2018b; Rajeswaran et al. 2017]. To encourage diversity during early training, recent work on transfer learning has also explored a multi-task paradigm where a model is pre-trained on many tasks before being transferred to a new target domain [Alet et al. 2018; Devin et al. 2017]. Some multi-task transfer learning solutions include policy distillation that seeks to “distill” knowledge from expert policies to a target policy [Parisotto et al. 2016; Rusu et al. 2016a]. Another approach with a similar goal is policy learning which learns a residual around given expert policies [Silver et al. 2019].

Meta learning has also gained popularity recently in computer vision and robotics, seeking to leverage past experiences obtained from many tasks to acquire a more generalizable and faster model that can be quickly adapted to new tasks [Andrychowicz et al. 2016; Ravi and Larochelle 2017]. The related formulations can be broadly classified into models that ingest a history of past experiences through recurrent architectures [Duan et al. 2016; Heess et al. 2015], model-agnostic meta-learning methods [Finn et al. 2017; Nichol et al. 2018], and approaches for meta-learning hyperparameters, loss functions, and task-dependent exploration strategies [Gupta et al. 2018; Houthoofd et al. 2018; Xu et al. 2018].

While some of the aforementioned approaches have shown great promise for agent control problems, in this paper, we propose an approach that can quickly adapt RL policies for physically simulated humanoids through fine control tuning as well as augmentation injected in the latent space, loosely inspired by recent findings in image diffusion [Hu et al. 2021; Mou et al. 2023; Zhang and Agrawala 2023]. In character animation, related work has focused on motion style transfer tasks for *kinematic* characters [Aberman et al. 2020; Mason et al. 2018] and the recent work of Starke et al. [2022] shows exciting results about how a well-learned latent space can aid motion synthesis. However, in physics-based character control tasks, there is still little investigation about the latent space representation of the

control policy obtained using reinforcement learning frameworks. We believe that AdaptNet provides a promising step in bridging that gap.

2.3 Latent Space Adaptation

We are inspired by research in image and 3D model generation that shows it is possible to control the synthesis process to generate targeted artifacts through purposeful modification of the latent space [Abdal et al. 2019; Berthelot et al. 2017; Bojanowski et al. 2018; Epstein et al. 2022; Karras et al. 2020; Radford et al. 2016; Shen et al. 2020; Wu et al. 2016; Zhuang et al. 2021]. While we have seen related work in RL for character control, AdaptNet offers a unique approach to latent space adaptation, drawn from these adjacent works’ successes. Related works in physics-based character control, such as [Juravsky et al. 2022; Ling et al. 2020; Peng et al. 2019, 2022; Tessler et al. 2023; Won et al. 2021], explore using pre-trained latent space models to facilitate the training of a control policy. These methods intend to adapt the pre-trained multi-skill model for downstream tasks by controlling skill latent embeddings, focusing on reusing skills for motion generation. In contrast, our approach does not break down the latent space by task and character state and instead allows the policy to be adapted to heterogeneous tasks that require learning new (out-of-distribution) motions/skills. Further, previous methods discard the pre-trained latent encoder during adaptation and rely on re-training to obtain a new encoder. In contrast, our approach directly edits the latent space projected from the association of the task and character state via the pre-trained policy. To do this, we use a gated recurrent unit (GRU) [Chung et al. 2014] layer as the encoder and initialize it by duplicating the original encoder parameters. Next, a fully connected layer is applied after the GRU to ensure zero initialization and convert the encoded state to a latent *modification*. In sum, the training for our adaptation starts from modifying the pre-trained policy rather than from scratch, which benefits adaptation in comparison to previous work in sample efficiency and, at times, overall effectiveness.

3 ADAPTNET FRAMEWORK

An overview of the AdaptNet framework is shown in Figure 2. The GAN-style control framework (top), described below, produces an original (pre-trained) policy (bottom, left) while AdaptNet is used to adapt that pre-trained control policy to a new task controller (bottom, right). Notably, the adaptation process could involve changes to the reward function (e.g., motion stylization) or the state and dynamics model (e.g., character morphology and terrain adaptation). Components of the AdaptNet for policy adaptation are shown: a latent space injection component and an internal adaptation component. The latent space injection performs policy adaption by editing the latent space, which is conditioned on the pre-trained policy’s state as well as any additional state information, for example, for new tasks. This component is trained to cooperate with the pre-trained policy by generating offsets to the original latent space instead of trying to learn how to generate latent variables for new tasks from scratch during adaptation. This leads to an efficient state-action exploration that starts from the pre-trained policy, instead of complete random exploration. Internal adaptation further tunes the policy by

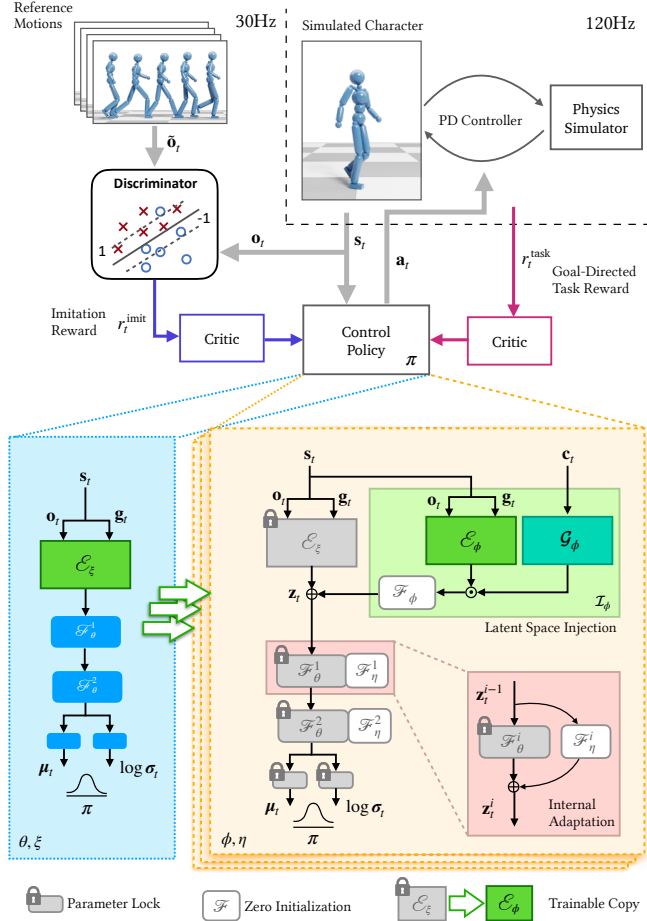


Fig. 2. Overview of our approach for adapting motor control policies for physics-based characters. Top: We model both pretraining and adapted tasks using a multi-critic reinforcement learning framework that balances the training of imitation and goal-directed control objectives. After a policy is trained, we can quickly adapt it to a new task using AdaptNet. Bottom: AdaptNet starts with a copy of the pre-trained policy network and modifies it through editing the latent space conditioned on the character's state and introducing optional adaptation modules for further finetuning.

adding a branch to each internal fully-connected layer in the policy network. This allows for more flexibility, enabling AdaptNet to shift away from the pre-trained policy and generate refinement through control actions that the pre-trained policy may not reach easily.

In our implementation, both the pre-trained policy and the adaptation are produced using a multi-objective learning framework [Xu et al. 2023] combining reinforcement learning with a GAN-like structure for effective policy learning that accounts for both motion imitation and goal-directed control (see Figure 2, top). During runtime, AdaptNet can be activated flexibly and dynamically allowing us to control the level of adaptation of the original control policy. The control policy $\pi(a|s_t)$ is a neural network taking the agent state s_t as input and outputting a probability distribution from which a control a_t can be drawn from the action space \mathcal{A} .

For physics-based character control tasks with dynamic goals, we consider $s_t := \{o_t, g_t\}$, where o_t denotes the current state of the character, e.g., joint or body link positions and velocities, and g_t is an optional task-related goal state or an encoding variable that indicates desired motion parameters, such as target speed and direction, end-effector positions, motion style, etc. The action vector a_t is the target posture fed to a PD servo through which the simulated character is controlled at a higher frequency. As shown in Figure 2, a_t is expressed as a multivariate Gaussian distribution.

Under the framework of reinforcement learning, our goal is to find the policy π that maximizes the discounted cumulative reward:

$$J = \mathbb{E}_{\tau \sim p(\tau|\pi)} \left[\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) \right], \quad (1)$$

where $p(\tau|\pi) = p(s_0) \prod_{t=0}^{H-1} p(s_{t+1}|s_t, a_t) \pi(a_t|s_t)$ is the state-action visitation distribution for the trajectory $\tau = \{s_t, a_t\}$ over a horizon of H time steps, $\gamma \in [0, 1]$ denotes the discount factor, and $r(\cdot)$ is the reward received at a given time step and $p(\cdot)$ is the state-transition probability of the underlying Markov decision process. In our domain, when the character faces a new task, $p(\cdot)$ and/or $r(\cdot)$ may change. AdaptNet seeks to efficiently modify π and adapt it to the new task by editing the latent space and finetuning the policy.

4 POLICY ADAPTATION USING LATENT SPACE INJECTION

If we consider the first layer, or first several layers, in the policy network π as an encoder to embed the state s_t into a latent space \mathcal{Z} , the control policy can be rewritten as

$$\pi_{\theta}(a_t | \mathcal{E}_{\xi}(s_t)), \quad (2)$$

where \mathcal{E}_{ξ} is the encoding layers with parameters ξ , θ are the parameters for the layers in the policy network that follow the encoder, and (θ, ξ) denote the weights of π . In this formulation, the policy network π_{θ} decides the projection from the latent $z_t = \mathcal{E}_{\xi}(s_t)$ into the action space \mathcal{A} . Assuming that π_{θ} is optimized by a typical on-policy policy gradient algorithm, the optimization objective with the introduction of the latent becomes

$$\max_{\theta, \xi} \mathbb{E}_t [A(s_t, a_t) \log \pi_{\theta}(a_t | z_t; \xi)], \quad (3)$$

where $A(\cdot)$ provides an advantage function estimation based on the received rewards $\{r_k\}_{k \geq t}$ during the interaction with the environment and represents how good an action sample a_t is given the conditional state s_t .

Given the generalization of neural networks, the latent space \mathcal{Z} can be considered as a superset covering all the possible latent states, which could lie outside of the domain that π_{θ} can reach during its training.

Based on this observation, when π_{θ} needs to be adapted to a new task, we propose to edit $z_t = \mathcal{E}_{\xi}(s_t) \subset \mathcal{Z}$ instead of discarding the original encoder \mathcal{E}_{ξ} and training a new one from scratch. The intuition is that for similar tasks, adjusting the current encoder provides better efficiency, allowing the desired control policy to be learned by a modified projection function from s_t to z_t .

Our approach manipulates the full latent space projected from both the character state o_t and the goal state g_t . Specifically, as

shown in Figure 2, we perform latent space injection by introducing a new conditional encoder \mathcal{I}_ϕ with parameters ϕ after the first encoding layer, where the character state \mathbf{o}_t and the goal state \mathbf{g}_t are concatenated to generate \mathcal{E}_ξ . This latent space is modified via

$$\mathbf{z}_t = \mathcal{E}_\xi(\mathbf{s}_t) + \mathcal{I}_\phi(\mathbf{s}_t, \mathbf{c}_t), \quad (4)$$

where \mathbf{c}_t is an additional control input for the new task which could be optional. The injector module \mathcal{I}_ϕ is

$$\mathcal{I}_\phi(\mathbf{s}_t, \mathbf{c}_t) = \mathcal{F}_\phi(\text{CONCAT}(\mathcal{E}_\phi(\mathbf{s}_t), \mathcal{G}_\phi(\mathbf{c}_t))), \quad (5)$$

where \mathcal{G}_ϕ is an optional module to process the additional control input \mathbf{c}_t , \mathcal{E}_ϕ is a state encoder that has exactly the same structure as the original encoder \mathcal{E}_ξ , and \mathcal{F}_ϕ is a final embedding module, which can be a fully-connected layer or a stack of multiple fully-connected layers.

During retraining for adaptation, we perform policy optimization as in Eq. 3, but only optimize the new parameters ϕ while keeping the parameters θ and ξ fixed:

$$\max_{\phi} \mathbb{E}_t [A(\mathbf{s}_t, \mathbf{a}_t) \log \pi_\theta(\mathbf{a}_t | \mathcal{E}_\xi(\mathbf{s}_t) + \mathcal{I}_\phi(\mathbf{s}_t, \mathbf{c}_t))]. \quad (6)$$

We begin with copying the original encoder parameters ξ into the new encoder \mathcal{E}_ϕ and initializing the last fully-connected layer inside \mathcal{F}_ϕ with zero weight and bias. In this way, the new encoder \mathcal{E}_ϕ is optimized by finetuning a set of parameters that are already optimized for state feature extraction during pre-training. The zero initialization of \mathcal{F}_ϕ lets the control policy give exactly the same action output as the original pre-trained one, i.e., $\pi_\theta(\mathbf{a}_t | \mathcal{E}_\xi(\mathbf{s}_t))$, at the beginning of re-training. It guides the adaptation to start from the state-action trajectory generated by the original policy rather than from a completely random exploration.

We refer to Figure 2 for the default implementation of AdaptNet, where the latent space injection is performed right after the concatenation of \mathbf{o}_t and \mathbf{g}_t . We denote this latent space as \mathcal{Z}^0 , and the following ones after each fully-connected layer but before the final action layer as \mathcal{Z}^i where $i = 1, 2, \dots$. Empirically, we note that it is more challenging to perform optimization when the injection occurs at a deeper layer in the policy network, leading typically to unstable training and low-fidelity controllers. An extreme case is to perform injection directly at the action space, which makes the whole system similar to directly finetuning the pre-trained policy network. We refer to Section 9 for related sensitivity analysis on introducing latent space injection at different network layers and for comparisons with directly finetuning a pre-trained policy network for new tasks.

During runtime, we can further introduce an extra scaling coefficient to the injection term in Eq. 4. Since our approach does not change the original encoder \mathcal{E}_ξ as well as the policy π_θ , the scale coefficient allows us to turn the injection on and off, or control the transition from the original policy to the fully adapted one. In such a way, we can perform motion style or behavior transitions (e.g., walk to skip) by interpolation in the latent space, as we will show in Section 8.1.

5 INTERNAL ADAPTATION FOR CONTROL LAYERS

The latent space injection component of AdaptNet edits the latent space based on the input state and further allows us to introduce

additional control input for new tasks. However, the expressive ability of the action policy is still constrained by the pre-trained layers after the state encoder in the policy network, i.e., π_θ . While utilizing the pre-trained π_θ for fast adaptation to new tasks, we introduce an internal adaptation component through which we can finetune π_θ , overcoming the bias it introduces and allowing for more flexibility in the types of generated controls compared to the ones obtained from the original training domain. The goal of the finetuning is to find a *small* increment $\Delta \mathbf{z}_t^i$ to the original latent \mathbf{z}_t^i in each latent space \mathcal{Z}^i , $i > 1$, to help optimize the objective function in Eq. 6 during adaptation training, but without changing the π_θ too much to avoid drifting too far away from the pre-trained policy and being stuck at overfitting during adaptation. To do so, we add a branch to each fully-connected layer between two latent spaces. As shown in the red block of Figure 2, the corresponding latent is generated as:

$$\mathbf{z}_t^i = \mathcal{F}_\theta^i(\mathbf{z}_t^{i-1}) + \mathcal{F}_\eta^i(\mathbf{z}_t^{i-1}). \quad (7)$$

Here, \mathcal{F}_θ^i denotes the fully-connected layer between the latent space \mathcal{Z}^{i-1} and \mathcal{Z}^i in the policy network π_θ , and \mathcal{F}_η^i is the newly introduced adaptor that generates $\Delta \mathbf{z}_t^i$ and is modeled as a fully-connected layer in the added branch. The parameter η is defined as

$$\eta := \{\Delta \mathbf{W}_i, \Delta \mathbf{b}_i\}, \quad (8)$$

with $\Delta \mathbf{W}_i$ and $\Delta \mathbf{b}_i$ being the weight and bias parameters in \mathcal{F}_η^i respectively. Similarly to the embedding module \mathcal{F}_ϕ in the latent space injection component, \mathcal{F}_η^i is initialized as zero and will not influence the output of the policy network at the beginning of policy adaptation. We lock θ in \mathcal{F}_θ^i during adaptation training and introduce the parameter η into the optimization function in Eq. 6.

Our approach is different from directly finetuning π_θ . When directly finetuning π_θ , the gradient from \mathbf{z}_t^i with respect to \mathbf{z}_t^{i-1} is decided by the weight \mathbf{W}_i in the layer \mathcal{F}_θ^i , which may be highly biased and have relatively large or very small values given it was fully trained. Therefore, finetuning π_θ directly for new tasks may lead to unstable training compared to only finetuning the newly introduced parameter set η which is initialized with zero. Furthermore, we can easily apply regularization on $\Delta \mathbf{W}_i$ and $\Delta \mathbf{b}_i$ to prevent aggressive finetuning regardless of the value of the parameters \mathbf{W}_i and \mathbf{b}_i in the pre-trained layer \mathcal{F}_θ^i . This will limit the possible change that the internal adaptation can bring about in order to prevent overfitting. We can also introduce an extra scaling weight to control the adaptation level during runtime, as discussed in Section 4.

Our proposed internal adaptation component is similar to the approach of low-rank adaptation (LoRA) proposed by Hu et al. [2021]. The major difference is that instead of directly employing a fully-connected layer, LoRA decomposes the weight matrix $\Delta \mathbf{W}_i$ into two low-rank matrices, i.e., $\Delta \mathbf{W}_i = \mathbf{B}_i \mathbf{A}_i$, where, \mathbf{B}_i is a $|\mathcal{Z}^{i-1}|$ -by- r matrix, \mathbf{A}_i is a r -by- $|\mathcal{Z}^i|$ matrix, and $r \ll \min(|\mathcal{Z}^{i-1}|, |\mathcal{Z}^i|)$. In contrast, our approach can be considered a full-rank adaptation. LoRA has been demonstrated as an effective way to fine tune large language and image generation models, reducing the number of parameters that need to be optimized during model adaptation. However, as shown in Section 9.3, we found that LoRA does not work well for physics-based character control tasks. A possible reason is

ALGORITHM 1: Policy Adaptation using AdaptNet

-
- Obtain the policy π_θ and the state encoder \mathcal{E}_ξ by performing training to optimize Eq. 9 in a general or default environment setting.
- 1 Build up the latent space injection component \mathcal{I}_ϕ based on Eq. 5 and the internal adaptation component $\{\mathcal{F}_\eta^i\}$ based on the Eq. 7.
 - 2 Lock the parameters θ and ξ .
 - 3 Initialize \mathcal{E}_ϕ using the pre-trained parameter ξ .
 - 4 Initialize the last layer inside \mathcal{F}_ϕ and each \mathcal{F}_η^i using zero weight and bias.
 - 5 Adapt the policy for a new task by only optimizing the parameters ϕ and η using Eq. 12.
-

that the related policy networks are markedly smaller compared to large language and image generation models that may have more than 12K dimensions. The latent spaces of our policy network have a typical size of 512 or 1024 dimensions and may not exhibit the lower intrinsic ranks that larger models do [Aghajanyan et al. 2021; Li et al. 2018; Pope et al. 2021].

6 POLICY TRAINING

We use the multi-objective learning framework for physics-based character control proposed by Xu et al. [2023] to perform both the original (pre-)training and adaptation training. The framework leverages a multi-critic structure where the objectives of motion imitation and goal-directed control are considered independent tasks during policy updating. In Figure 2, for example, the imitation objective is associated with a critic network labeled in blue, and the goal-directed objective is associated with a critic in magenta. The advantage (cf. Eqs. 3, 6) with respect to each objective is estimated only by its associated reward and critic network. To ensure that the policy can be updated in a balanced way taking into account both the imitation and goal-directed control objectives, all estimated advantages are standardized independently before policy updating.

During pre-training, we seek to find a basic motor control policy $\pi_\theta(\mathbf{a}_t|\mathcal{E}_\xi(\mathbf{s}_t))$, which we can later adapt to new tasks. In this work, we focus on locomotion tasks, and thus π_θ involves two objectives: a motion imitation objective given a batch of reference motions of walking and running, and a goal-directed objective involving a given target direction and speed. Using the multi-objective learning framework, the optimization objective function during pretraining shown in Eq. 3 can be written as

$$\max_{\theta, \xi} \mathbb{E}_t \left[\left(\sum_k \omega_k \bar{A}_t^k \right) \log \pi_\theta(\mathbf{a}_t | \mathcal{E}_\xi(\mathbf{s}_t)) \right], \quad (9)$$

where \bar{A}_t^k is the standardization of the estimated advantage associated with the objective k and ω_k satisfies $\sum_k \omega_k = 1$ providing additional control to adjust the policy updating in a preferred manner when conflicts between multiple objectives occur.

We employ a GAN-like structure [Ho and Ermon 2016; Merel et al. 2017] that relies on an ensemble of discriminators [Xu and Karamouzas 2021] to evaluate the imitation performance and generate the corresponding reward signals for advantage estimation and policy updating. In particular, we take an ensemble of N discriminators and use a hinge loss [Lim and Ye 2017] with policy gradient [Gulrajani et al. 2017] for discriminator training, resulting

in the following loss function:

$$\min \frac{1}{N} \sum_{n=1}^N \left(\mathbb{E}_t [\max(0, 1 + D_n(\mathbf{o}_t))] + \mathbb{E}_t [\max(0, 1 - D_n(\tilde{\mathbf{o}}_t))] + \lambda^{\text{GP}} \mathbb{E}_t [(\|\nabla_{\tilde{\mathbf{o}}_t} D_n(\tilde{\mathbf{o}}_t)\|_2 - 1)^2] \right). \quad (10)$$

Here, D_n denotes a discriminator network, $\tilde{\mathbf{o}}_t = \alpha \mathbf{o}_t + (1 - \alpha) \tilde{\mathbf{o}}_t$ with $\alpha \sim \text{UNIFORM}(0, 1)$ and λ^{GP} is gradient penalty coefficient. The reward function to evaluate the imitation performance is defined as

$$r^{\text{imit}}(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) = \frac{1}{N} \sum_{n=1}^N \text{CLIP}(D_n(\mathbf{o}_t), -1, 1). \quad (11)$$

The reward for the goal-related task is computed heuristically. We refer to the appendix for the representation of the goal state \mathbf{g}_t and the definition of the goal-related task reward.

After obtaining π_θ and \mathcal{E}_ξ in pre-training, we introduce the proposed AdaptNet to perform policy adaptation for new tasks that are relative to but have different reward definitions and/or environment settings from the one in the pre-training phase. Before the adaptation training starts, we lock the parameters θ and ξ . We then initialize \mathcal{E}_ϕ inside the latent space injection component \mathcal{I}_ϕ using the weights ξ , and initialize with zero weight and bias the last layer of \mathcal{F}_ϕ inside \mathcal{I}_ϕ along with each fully-rank adaptor \mathcal{F}_η^i , $i > 0$. To stabilize the training, besides applying a common weight decay to the parameter set η (Eq. 7) via L2 regularization, we introduce an additional regularization on the latent injection generated by \mathcal{I}_ϕ . The adaptation training is still performed under the aforementioned multi-objective learning framework in the same way as the pre-training phase. The optimization objective for policy adaptation is

$$\max_{\phi, \eta} \mathbb{E}_t \left[\left(\sum_k \omega_k \bar{A}_t^k \right) \log \pi_\theta(\mathbf{a}_t | \mathcal{E}_\xi(\mathbf{s}_t) + \mathcal{I}_\phi(\mathbf{s}_t, \mathbf{c}_t); \eta) - \beta \|\mathcal{I}_\phi(\mathbf{s}_t, \mathbf{c}_t)\|_2 - \kappa \|\eta\|_2 \right], \quad (12)$$

where β and κ are regularization coefficients. In Section 10, we give a detailed analysis of the regularization on the latent space injection.

We refer to Algorithm 1 for the outline of the whole training process. Adaptation with the proposed AdaptNet can be done very quickly within 10-30 minutes for simple control tasks and up to 4 hours for challenging terrain adaptation tasks with new control input processed by an additional convolutional neural network \mathcal{G}_ϕ , as defined in Eq. 5.

7 EXPERIMENTAL SETUP

Our experiments were run using IsaacGym [Makoviychuk et al. 2021] with 512 environments running in parallel during training. The simulated character has 15 body links and 28 degrees of freedom, where the elbow and knee joints are implemented as 1-dimensional revolute joints, and the hands are fused with the forearms and uncontrollable. All simulations run at 120Hz with a normal PD controller employed as the low-level actuator to directly manipulate the simulated character, while the control policy runs at 30 Hz, as shown in Figure 2.

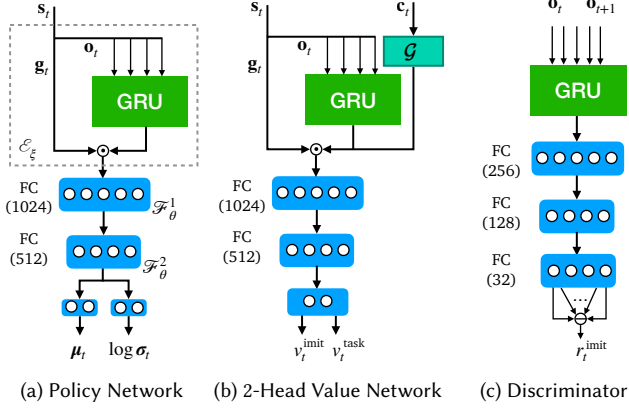


Fig. 3. Network structures. Here, \odot denotes the concatenation operator and \ominus denotes the average operator. The state encoder \mathcal{E}_ξ is shown in the dashed block. An optional control input encoding module \mathcal{G} is included if the additional control input c_t is provided during adaptation training.

We run policy optimization using PPO [Schulman et al. 2017] and update policy parameters using the Adam optimizer [Kingma and Ba 2017]. To encode the character’s state, we take the position, orientation, and velocities of all the body links related to the pelvis (root link) in the last four frames as the state representation \mathbf{o}_t and employ a gated recurrent unit (GRU) [Chung et al. 2014] with a 256-dimension hidden state to process this temporal state. For discriminator training, we take the character’s pose at five consecutive frames as the representation of $\{\mathbf{o}_t, \mathbf{o}_{t+1}\}$ to evaluate the policy’s imitation performance during the transition from timestep t to $t + 1$. We employ an ensemble of 32 discriminators and model it by a multi-head network, as shown in Figure 3. The critic network has a similar structure to the policy network, but with a 2-dimensional output for the value estimations to the imitation objective and goal-directed objective respectively. We refer to the appendix for the hyperparameters used for policy training and the representation of the goal state \mathbf{g}_t in the locomotion task.

Rewards for both task and imitation are employed during policy adaptation. To avoid bias from the pre-trained policy, we discard the discriminators for imitation from the original policy and new discriminators are trained from scratch. Intuitively, in tasks such as motion style transfer the original discriminator will not work well for the new given reference style and thus a new one is needed. Even for other adaptation tasks, we found utilizing old discriminators to be problematic, as the optimal action in the new task can dramatically change from the original in the context of how it employs the reference motion. Empirically, when we experimented with reusing the old discriminators, we found they introduce too much bias towards the old task. Finally, with training new discriminators for a new task, we also perform value estimation by re-training a new critic from scratch.

All our tests were run on machines with a V100 or A100 GPU. To achieve a good locomotion policy based on which we perform further adaptation, the pre-training took around 26 hours and consumed 4×10^8 training samples. The reference motions are around

Table 1. Reference motions for policy pre-training (top) and stylized motion learning (bottom).

Motion	Length	Description
Walk	334.07 s	normal walking motions for pre-training
Run	282.87 s	normal running motions for pre-training
Swaggering Walk	1.07 s	exaggerated walking with one arm akimbo
Goose Step	2.20 s	goose step with arms akimbo
Stomp Walk	1.23 s	walking while stomping on the ground
Kicking Walk	2.03 s	walking with leg kicking
Stoop	0.93 s	slow walking with body bent over
Jaunty Skip	1.60 s	skipping in a spirited manner
Sashay Walk	1.07 s	walking in a slightly exaggerated manner
Limp	1.90 s	slow walking with right leg hurt
Pace	1.70 s	slow walking with arms akimbo
Penguin Walk	0.77 s	moving with very small and steps
Strutting Walk	1.40 s	walking with shoulder moving aggressively
Joyful Walk	1.20 s	strut walking rhythmically

300 seconds long including normal walking and running motions with turning poses and various speeds (cf. Table 1, top). All the reference motions used during pre-training and adaptation training are recorded at 30 Hz and extracted from the publicly available dataset LAFAN1 [Harvey et al. 2020].

8 APPLICATIONS OF ADAPTNET

In this section, we apply the AdaptNet technique to demonstrate the success and efficiency of learning new physics-based controllers through adaptation. Our experiments use two pre-trained locomotion policies (walking and running) that account for two objectives: motion imitation based on a batch of walking or running reference motions, respectively, and a goal objective as defined by a target direction of motion and speed. We adapt the pre-trained policies to a range of new tasks, highlighting applications of AdaptNet to style transfer, character morphology changes and adaptation to different terrains. Figure 1 shows snapshots from different outcomes. Please refer to the supplementary video for related animation results.

8.1 Motion Style Transfer and Interpolation

We consider a variety of motion style transfer tasks where a pre-trained walking locomotion policy is adapted to a particular style. Note, this is not a simple motion imitation task, since all the style reference motions are very short (see Table 1, bottom), containing only one or two gait cycles. It is therefore impossible to train an equivalent locomotion policy that supports goal-directed steering using the target reference motion. Instead, the nature of this test is few-shot learning, where AdaptNet is expected to effectively learn how to perform locomotion in the style provided by the small duration of the style example in the new reference, while relying on the pre-trained policy to perform turning and goal-directed steering. Figure 5 depicts related qualitative results. AdaptNet can effectively learn how to do goal-directed turning in the provided style. Further, adaptation training can be done very quickly, within 10-30 minutes, in contrast to the original that we obtained during pre-training took about one day for training. We refer to the supplementary video for

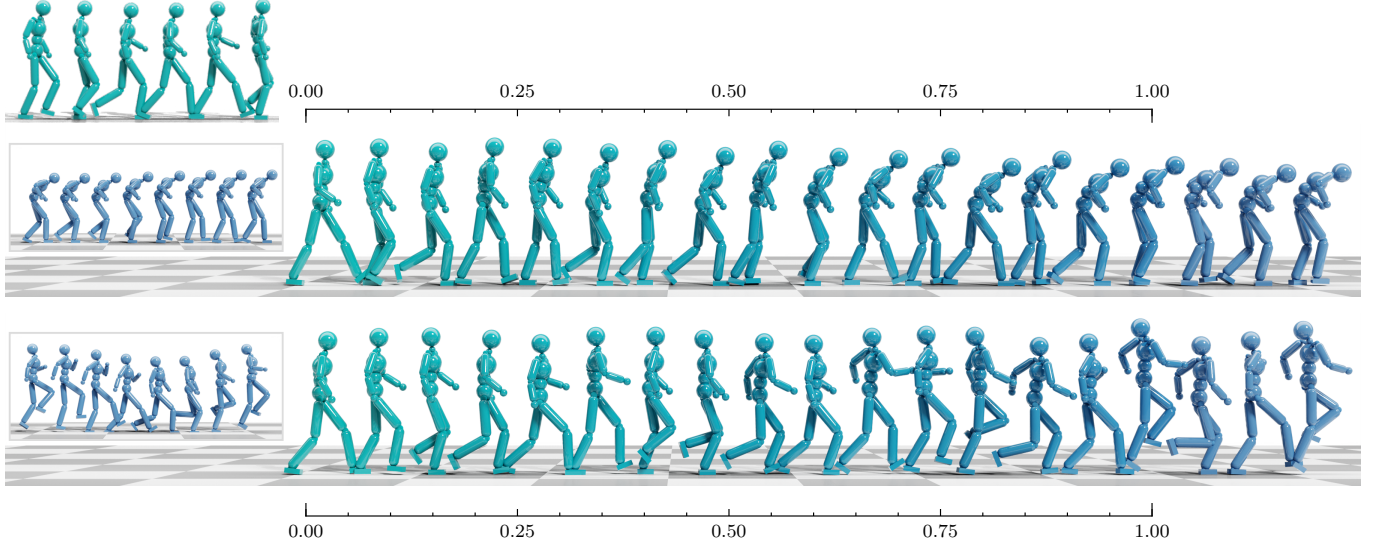


Fig. 4. Motion interpolation between walking (pre-trained policy) shown at the top-left corner and different stylized motions by controlling the adaptation level of the associated AdaptNet model (cf. Eq. 13). Snapshots on the left show the learned stylized motions of *Stoop walking* and *Jaunty Skip*. When $\alpha = 0$, the character is controlled only by the original walking policy. When $\alpha = 1$, the character is controlled with a full injection of AdaptNet.



Fig. 5. Example of motion style transfer learning with goal-steering navigation using AdaptNet. Green arrow indicates the dynamically generated target directions for locomotion control.

animation results, and Section 9 for comparing AdaptNet to learning stylized locomotion from scratch.

As discussed in Sections 4 and 5, we can perform motion interpolation in the latent space by introducing a scale variable to control the adaptation level. This process can be described by modifying Eqs. 4 and 7 as

$$\begin{aligned} \mathbf{z}_t^0 &= \mathcal{E}_{\xi}(\mathbf{s}_t) + \alpha \mathcal{I}_{\phi'}(\mathbf{s}_t, \mathbf{c}_t), \\ \mathbf{z}_t^i &= \mathcal{F}_{\theta}^i(\mathbf{z}_t^{i-1}) + \alpha \mathcal{F}_{\eta}^i(\mathbf{z}_t^{i-1}), \end{aligned} \quad (13)$$

where $\alpha \in [0, 1]$ is the introduced scale variable. In Figure 4, we show interpolation results. As shown in the figure, we can achieve motions with different style intensity, which can transition between the base walking motion and the stylized ones in a smooth manner.

We can further extend Eq. 13 to perform interpolation between any two AdaptNet models via

$$\begin{aligned} \mathbf{z}_t^0 &= \mathcal{E}_{\xi}(\mathbf{s}_t) + \alpha \mathcal{I}_{\phi'}(\mathbf{s}_t, \mathbf{c}_t) + (1 - \alpha) \mathcal{I}_{\phi''}(\mathbf{s}_t, \mathbf{c}_t), \\ \mathbf{z}_t^i &= \mathcal{F}_{\theta}^i(\mathbf{z}_t^{i-1}) + \alpha \mathcal{F}_{\eta}^i(\mathbf{z}_t^{i-1}) + (1 - \alpha) \mathcal{F}_{\eta''}^i(\mathbf{z}_t^{i-1}), \end{aligned} \quad (14)$$

where the parameters ϕ' and η' are from one AdaptNet model and ϕ'' and η'' are from the other one. Such an interpolation scheme can be regarded as applying two independently trained AdaptNet models simultaneously on the same, pre-trained policy, with an example shown in Figure 6.

The above interpolation results demonstrate that during adaptation training, AdaptNet can effectively learn structured information about the latent space with respect to the desired motion styles. We refer to Section 10 for more details on controlling the latent space and related visualizations, along with an analysis of the training difficulty (time consumption) when learning different styles.

8.2 Morphological Adaptation

We consider two kinds of morphological changes: body shape and joint lock. Due to physical constraints, morphological changes in the character model will cause the same action \mathbf{a}_t to lead to different resulting states compared to the ones observed in the pre-training phase. Without adaptation, the pre-trained policy does not perform well if it's even able to keep the character balanced, especially when the lower body is modified.

We tested eight body-shape variants of the original character model, as shown in Figure 7. In the *LongBody* variant, we extend the abdomen length by 50%, while the *BigBody* variant increased the torso size by 50%. The latter leads to an increase in the torso mass of over 300%. In *LongUpperArms* and *LongLowerArms* variants, the length of the upper and lower arms are extended by 25% respectively, while in *AsymmetricUpperArms*, we increase the length of the right upper arm but decrease the length of the left upper arm. In the *LongThighs* and *LongShins* variants, the length of the upper and lower legs are extended by 50% respectively, the latter akin to a human walking on stilts. In the model of *SuperLongLegs*, both the

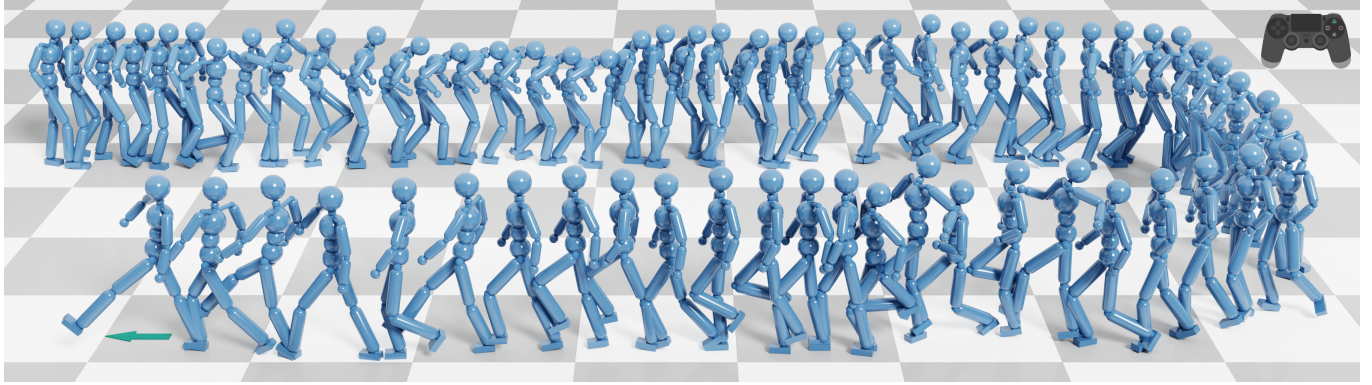


Fig. 6. Motion interpolation in the latent space by activating and switching between multiple AdaptNet models to let the character perform style transition interactively during goal-steering navigation.

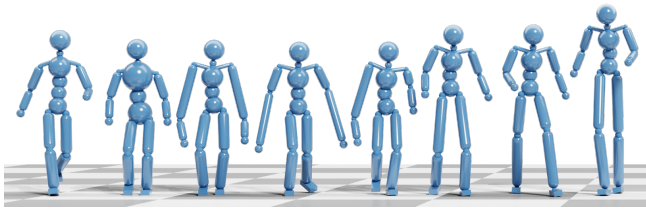


Fig. 7. Character models with body shape variants. From left to right: *LongBody*, *BigBody*, *LongUpperArms*, *LongLowerArms*, *AsymmetricUpperArms*, *LongThighs*, *LongShins*, and *SuperLongLegs*.

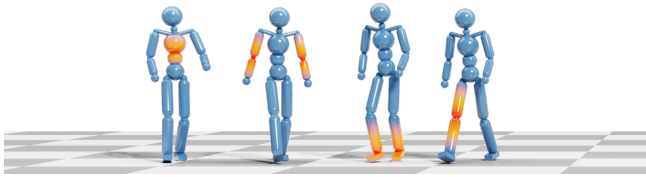


Fig. 8. Character models with joints being locked. From left to right, the locked joints are abdomen, elbows, ankles, and right knee respectively (shown in red). Corresponding body parts between a locked joint are highlighted in orange.

thighs and shins are extended resulting in a character that is over 2 m tall.

We also experimented with different configurations, as shown in Figure 8, where some of the joints (in orange) are ‘locked’. The locked joints are removed from the character model such that the linked body parts are fused together. This reduces the number of dimensions of the action space. To make the pre-trained policy compatible with the new action space, we simply prune the weight and bias matrices of the last layers in the policy network and remove the output neurons corresponding to the locked joints.

Even though the pre-trained policy would not completely lose control of the character when the torso or arms are modified, the character still loses balance quite often. As more challenging examples, the morphological changes in the lower body parts and joints leave the pre-trained policy unable to control the character without falling.

For example, when the knee joint is locked, the policy needs to adjust the output of the hip and ankle in order to compensate for the ‘disability’ of the knee. This requirement leaves the pre-trained policy incapable of suitably controlling the modified character model.

During adaptation, we did not do any retargeting to generate new reference motions for AdaptNet to learn. Instead, we simply modify the character’s model while relying on the reference motions used to pre-train the original policy, retargeted to the character model without any morphological changes. We found it takes 15-30 minutes to finish the adaptation training depending on the difficulty of the morphology change task. The character controlled by the AdaptNet policy can maintain its balance and walk or run without falling down. An interesting observation is that in order to match the provided height of the root link (pelvis) in the reference motions, the AdaptNet policy will control the character to walk or run in a crouch with the body at a relatively low position compared to the leg length. We show some representative results in Figure 9, and refer to the supplementary video for animations.

8.3 Terrain Adaptation

Next we discuss policy adaptation for character locomotion on low friction and rough terrains as well as obstacle-filled scenes that require extra control input.

8.3.1 Friction Adaptation. To simulate an icy surface, we significantly reduce the ground friction. In particular, we decrease the friction coefficient from 1 to 0.15 for walking and to 0.35 for running. Figure 10 compares results obtained for the running policy with and without using AdaptNet. Note, AdaptNet can effectively control the character to change its moving direction by sliding on its feet, as shown in the left example of the figure. In addition, using AdaptNet,



Fig. 9. Adapting the locomotion policy of running to characters with different body shapes and locked joints.



Fig. 10. Comparison of characters controlled with and without AdaptNet running on an ice floor with very low friction. Left: character controlled with AdaptNet slides and skids on the ice ground while running. Right: character without AdaptNet slips down.

the character lowers its center of mass and takes quick steps to maintain its balance. In contrast, with the original policy, the character cannot run on the icy ground without falling down. For walking, the AdaptNet controller is more cautious with the character preferring to stop and change its direction in place. Without using AdaptNet, the character tends to turn around with a bigger radius, but not slow down. This demonstrates the ability of AdaptNet to change the behavior provided by the original policy to make it better suited to new environmental settings.

8.3.2 Terrain Adaptation with Additional Control Input. To test AdaptNet with extra control input, we designed several experiments where the character is asked to do goal-steering navigation in challenging environments with procedurally generated terrains. A local heightmap is provided as the additional control input c_t through which the character is expected to adjust its motions to prevent falling down during walking. The heightmap is extracted locally based on the character's root position and aligned with the orientation of the root, with a left and right horizon of 1.7 m, backward horizon of 1 m and forward horizon of 2.4 m. To process the heightmap c_t , we introduce a convolutional neural network (CNN) as the encoding module \mathcal{G}_ϕ (see Eq. 5) for AdaptNet. We refer to the appendix for the network structure of the CNN. An extra map encoding module having the same structure with \mathcal{G}_ϕ is added to the critic network for value estimation during adaptation. We show representative examples of our tested terrains in Figure 11 and note the appendix also gives more detail on terrain.

We refer to the companion video for the navigation performance of the character when walking on the designed terrains after adaptation training. Even in terrains where the height changes smoothly, the character teeters under the control of the pre-trained policy and a minor change in the terrain slope is enough to make the character stumble. After adaptation training, AdaptNet can enable the character to smoothly walk and turn on the uneven terrains

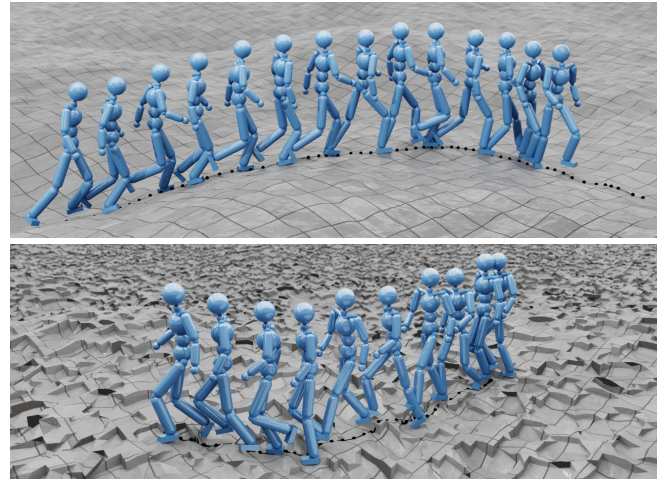


Fig. 11. Character controlled with AdaptNet navigates in the environment with procedurally generated terrains.

without falling. Besides being able to step over low-height obstacles, the AdaptNet character exhibits intelligent local decision making, trying not to step on the edge of the rocks on the rough terrain and avoids overly rugged paths by altering its moving trajectory to an easy-to-follow one.

To further demonstrate the ability of AdaptNet to perform local path planning, we designed a more challenging environment with uncrossable obstacles randomly placed on the ground. We qualitatively show the results in Figure 12. As seen in the figure, the character controlled with AdaptNet (blue) can successfully walk around the obstacles. Without accounting for collisions, the character controlled solely by the initially trained policy (green) crosses through the regions where obstacles are placed.

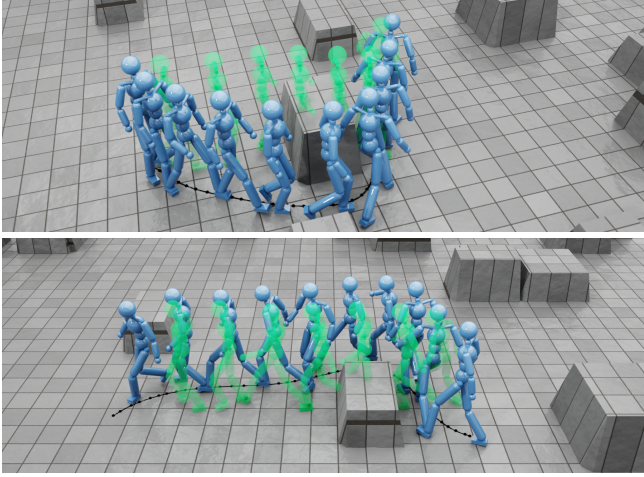


Fig. 12. Local collision avoidance in an obstacle-filled environment using AdaptNet. Green characters show the movement trajectory generated by the original walking policy without AdaptNet.

Unsurprisingly, the introduction of the CNN (detailed in Appendix B) increases the time needed to perform policy optimization iterations in the training for rough terrains. Still, for the easier terrains, training can be done within 1.5 hours. The more rugged terrain took around 4 hours for training. Finally, it took around 22 hours to train adaptation for the local obstacle avoidance test case. We note that this is still less time than is needed for training the original flat-ground locomotion policy from scratch (26 hours).

8.4 Perturbation Adaptation

In a final experimental foray, we investigate AdaptNet’s ability to improve the handling of perturbations. Although the original policy can handle small perturbations, the character will still fall under larger impulses. In order to achieve more robust control, we adapt the control policy’s ability to maintain balance in the presence of large disturbances. We begin with pre-trained policies for target-directed locomotion for walking and running. During the training process, we randomly apply perturbations (1000 N, lasting for 0.2 seconds) in different directions on the character’s torso. With adaptation training of around 5 hours, the character is able to stay balanced against comparable impulses following training for both running and walking tasks. In contrast, the original controls are not able to handle such perturbations repeatably and they often lead to the characters falling over. Furthermore, we also observe that AdaptNet control adjusts the character’s footsteps to recover balance when the character is highly out of balance due to perturbations. A comparison of the original policy and our results can be seen in the supplementary video.

9 ABLATION STUDIES

In this Section, we compare the performance of AdaptNet to different baselines along with performing sensitivity analysis on the two components of the proposed AdaptNet technique.

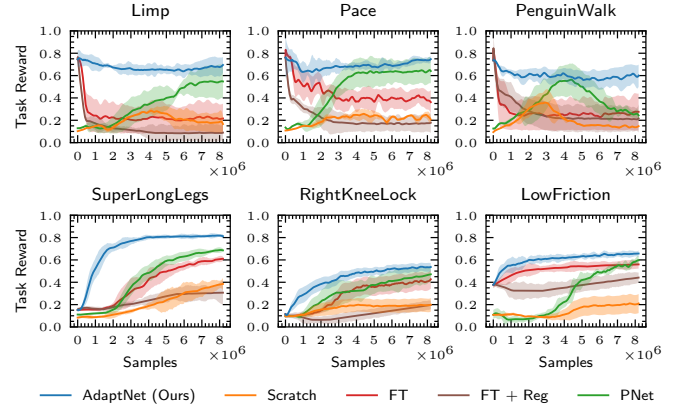


Fig. 13. Learning performance of our adaptation scheme using AdaptNet, training from scratch for each task (Scratch), using a progressive network (PNet), and adaptation via directly finetuning the pre-trained policy (FT) and finetuning with regularization (FT + Reg). Colored regions denote mean values \pm a standard deviation based on 5 trials. The top row consists of motion style transfer tasks, while the bottom row focuses on morphological and terrain adaptation tasks.

9.1 Baseline Comparisons

We consider the following baselines: *Scratch* where a new policy is trained from scratch on a given task; *FT* where we directly finetune the pre-trained policy network to the newly given task; *FT + Reg* where we apply regularization on the weights of the policy network during finetuning; and *PNet* where policy adaptation is performed using a progressive neural network approach [Rusu et al. 2016b]. Figure 13 compares the learning curves for the goal-task performance between the baselines and AdaptNet on three style-transfer tasks (top row) and three adaptation tasks (bottom row), two involving changes in the character’s morphology and one for lowered ground friction. For fair comparison, we employ the same training setup for all baselines, where the reward function of the new policy accounts for both a task objective and an imitation objective using an automatic weighting scheme [Xu et al. 2023]. In the motion style transfer experiments, the imitation term is computed using a new discriminator that takes only the stylized motions as the reference similar to Section 8.1.

As can be seen from the learning curves in Figure 13, *Scratch* fails to attain the desired goals in the considered benchmarks, achieving a very low goal task reward within the given budget of 8M training samples. *FT* can effectively modify the locomotion policy in the bottom three tasks where the character’s morphology or environmental friction changes. However, in the motion style transfer tasks, the reward curve of *FT* noticeably drops after the training begins as *FT* overfits the imitation of the newly provided stylized reference motion and ignores the goal direction signal. In contrast, AdaptNet provides a stable task reward curve during the adaptation training with the character being able to imitate the newly provided style without forgetting the previously learned locomotion behaviors as seen in Figure 14. The above findings are in line with previous works [Peng et al. 2019; Rusu et al. 2016b] that have shown finetuning to be efficient when the parameters of a pre-trained model need

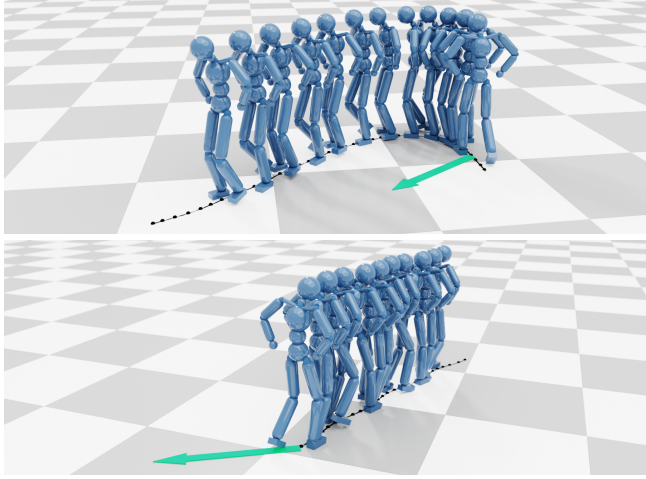


Fig. 14. Top: AdaptNet successfully controls the character to turn during walking in *Pace* style. Bottom: The character controlled by FT policy keeps imitating the reference motion to pace straightly, and fails to turn due to overfitting. Green arrows indicate the dynamically generated target directions for locomotion control.

to be slightly adjusted to a new target domain. However, *FT* can be susceptible to catastrophic forgetting when the imitation objective is significantly changed, as in the motion style transfer tasks. *FT + Reg* leads to poor training and low-fidelity controllers in all tasks. While, in theory, adding regularization can improve the navigation performance, in practice, it is hard to regulate the weights during finetuning due to the presence of both significant large and small weights in the pre-trained policy.

PNet shares similarities with AdaptNet as both approaches add new weights to the original policy network and freeze the old weights during transfer learning. However, despite these similarities, the architectures of the two approaches are significantly different. AdaptNet uses a residual structure that supports merging, resulting in a single policy network which allows forward propagation in one pass during inference. In contrast, *PNet* does not support merging and requires the original network to be present and run first to compute the values of the hidden neurons in the added network. This adds significant complexity and memory overhead, with the network structure becoming larger and slower. Importantly, during training, the added network in *PNet* cannot start from zero as compared to AdaptNet. In essence, the zero initialization in AdaptNet allows us to guide the adaptation starting from the original policy. This is clear in the style-transfer tasks, where AdaptNet begins training with a much higher reward than *PPNet* due to the locomotion ability provided by the original policy. Despite its competitive final performance in several of the adaptation tasks, *PNet* is sample inefficient. Finally, we note that it can lead to forgetting the prior knowledge provided by the pre-trained policy as the added network can significantly change the output of the whole model in some cases. This can be seen in the *Penguin Walk* task where the navigation performance drops after 5M samples.

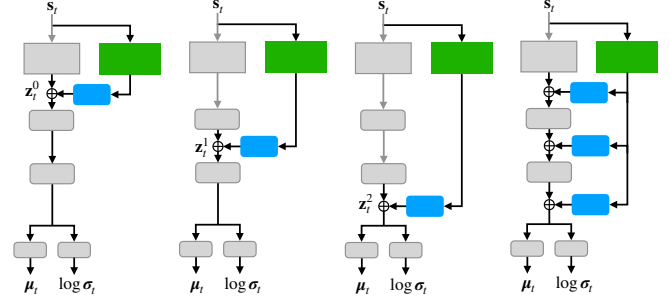


Fig. 15. Injection at different latent spaces. Gray blocks represent the original policy network locked during adaptation. Green blocks are the state encoder \mathcal{E}_ϕ , and blue ones are \mathcal{F}_ϕ . For left tor right, the manipulated latent spaces are \mathcal{Z}^0 (the default implementation of AdaptNet), \mathcal{Z}^1 , \mathcal{Z}^2 and $\mathcal{Z}^{0:2}$ respectively. We ignore \mathcal{G}_ϕ , given that there is no extra control input in the tested examples here.

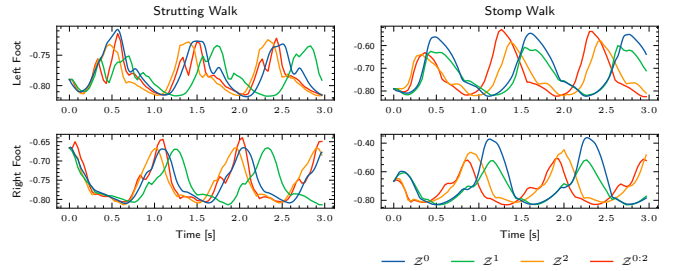


Fig. 16. Foot height (in meters) relative to the root when performing adaptation for motion style transfer tasks with different latent spaces being injected. Injection at \mathcal{Z}^0 (blue) leads to the smoothest and most repeatable stepping motions.

Overall, AdaptNet consistently outperforms all four baselines in terms of final performance and sample efficiency. In terms of memory efficiency, *Scratch* and *FT* do not add any overhead. AdaptNet introduces additional parameters, but since the original network is frozen, the number of trainable parameters is still at the same scale with the original neural network when no conditional input, i.e., c_t and \mathcal{G}_ϕ , is needed. While the total number of parameters increases, the effective number of parameters is the same as the original policy because AdaptNet can be merged into the original network. In contrast, *PNet* requires both networks to be present and effectively doubles the number of parameters.

9.2 Latent Space Injection

Our default implementation performs injection on the latent space \mathcal{Z}^0 right after the goal state \mathbf{g}_t and character state \mathbf{o}_t are encoded and concatenated together. Here, we test the application of the injection module to other latent spaces after \mathcal{Z}^0 but before reaching the action space, along with applying injection on all possible latent spaces simultaneously. To solely study the performance of latent space injection, we also remove the full-rank adaptation modules for these tests. The tested network structures are shown in Figure 15.

To explore how the injection schemes perform differently in generating new policies, we run tests on several motion style transfer tasks. During our experiments, we observe qualitatively that injection at the lower space \mathcal{Z}^2 or at all the latent spaces $\mathcal{Z}^{0:2}$, which also includes the lower one, can easily produce jerky motions with stiff movements of the torso and legs. It can also lead to failures in training where the character falls repeatedly after a few training iterations. In Figure 16, we plot the trajectory of the foot height in two of our tested cases. While injection at \mathcal{Z}^0 (blue) leads to a smooth repeatable trajectory, the curves become more irregular as the injected latent space changes from \mathcal{Z}^1 (green) to \mathcal{Z}^2 (orange) and then to $\mathcal{Z}^{0:2}$ (red). We also see some sharp jumps in the curves of \mathcal{Z}^2 and $\mathcal{Z}^{0:2}$, which represent fast motion transitions. We refer to the supplementary video for the animation results including examples where injection at \mathcal{Z}^2 and $\mathcal{Z}^{0:2}$ fails.

Overall, our tests show that as the chosen target latent space is closer to the action space, it becomes more difficult for AdaptNet to generate desired motions, with \mathcal{Z}^0 both intuitively and empirically giving the best results. This observation is in agreement with recent work in image synthesis where the target space for manipulation is usually chosen nearer to the input of the generator rather than near the final output [Abdal et al. 2019; Karras et al. 2020; Zhuang et al. 2021]. In terms of the network structure in our implementation, the input state $\mathbf{s}_t \in \mathbb{R}^{784}$ is encoded into the first latent space $\mathcal{Z}^0 \in \mathbb{R}^{260}$ and then projected to $\mathcal{Z}^1 \in \mathbb{R}^{1024}$. The whole network, therefore, can be regarded as an encoder-decoder structure where the bottleneck is at \mathcal{Z}^0 . As we will show in Section 10, \mathcal{Z}^0 is well-structured which makes it amenable to manipulation for motion generation.

9.3 Comparison of Adaptation Methods

We quantitatively evaluate the imitation performance of AdaptNet with other adaptation approaches, including alternate methods with and without using its internal adaptation component. As in prior work [Harada et al. 2004; Peng et al. 2021; Tang et al. 2008; Xu and Karamouzas 2021], we measure the imitation error via:

$$e_t = \frac{1}{N_{\text{link}}} \sum_{l=1}^{N_{\text{link}}} \|p_l - \tilde{p}_l\|, \quad (15)$$

where N_{link} is the total number of body links, $p_l \in \mathbb{R}^3$ is the position of the body link l in the world space at the time step t , and \tilde{p}_l is the body link's position in the reference motion. The evaluation results are shown in Table 2.

We find our proposed approach to combine latent space adaptation (LSA) and internal adaptation (IA) results in the best performance. While the results in Table 2 imply LSA alone is sufficient in many cases, IA appears to help most in the difficult motion style transfer tasks, e.g., *Goose Step*, *Jaunty Skip* and *Joyful Walk*, where the stylized motions are relatively far away from the pre-trained walking motions. In these tasks, adding IA improves the visual quality as well as motion smoothness, foot height, and gait frequency as shown in the supplementary video. It is important to note, however, that IA alone produces subpar performance. In addition, it cannot account for the additional control input needed in other adaptation tasks such as terrain adaptation. Further, even when no additional

Table 2. Imitation error during motion style transfer with different adaptation components. Values are reported in meters in the format of mean \pm std.

Motion	AdaptNet (LSA+IA)	LSA +LoRA-64	LSA	LSA w/o \mathcal{E}_ϕ	IA
Swagging	0.05 \pm 0.02	0.05 \pm 0.02	0.06 \pm 0.02	0.11 \pm 0.03	0.11 \pm 0.03
Goose Step	0.11 \pm 0.08	0.18 \pm 0.08	0.21 \pm 0.12	0.35 \pm 0.11	0.36 \pm 0.11
Stomp	0.08 \pm 0.04	0.10 \pm 0.05	0.11 \pm 0.06	0.26 \pm 0.07	0.27 \pm 0.08
Kicking	0.08 \pm 0.03	0.08 \pm 0.03	0.09 \pm 0.05	0.20 \pm 0.07	0.21 \pm 0.07
Stoop	0.07 \pm 0.02	0.07 \pm 0.02	0.07 \pm 0.02	0.14 \pm 0.03	0.13 \pm 0.03
Jaunty Skip	0.16 \pm 0.09	0.22 \pm 0.10	0.25 \pm 0.12	0.56 \pm 0.18	0.61 \pm 0.21
Sashay	0.06 \pm 0.03	0.06 \pm 0.03	0.06 \pm 0.03	0.09 \pm 0.03	0.09 \pm 0.04
Limp	0.10 \pm 0.07	0.10 \pm 0.07	0.12 \pm 0.07	0.22 \pm 0.09	0.29 \pm 0.11
Pace	0.09 \pm 0.03	0.10 \pm 0.03	0.10 \pm 0.03	0.14 \pm 0.03	0.13 \pm 0.03
Penguin	0.11 \pm 0.04	0.13 \pm 0.05	0.15 \pm 0.05	0.31 \pm 0.09	0.38 \pm 0.13
Strutting	0.09 \pm 0.03	0.10 \pm 0.05	0.12 \pm 0.06	0.23 \pm 0.06	0.27 \pm 0.06
Joyful	0.17 \pm 0.07	0.22 \pm 0.09	0.28 \pm 0.12	0.54 \pm 0.22	0.59 \pm 0.24

input is needed, IA components cannot be applied for modification of the state encoder as we cannot initialize the GRU layer of the encoder to zero. Latent modification is a distinctive feature of LSA, rendering Eqs. 4 and 7 unsuitable to be merged into the same formulation. To highlight the importance of the state encoder module \mathcal{E}_ϕ in LSA, we consider an additional ablation where we remove \mathcal{E}_ϕ and connect the output of \mathcal{E}_ξ directly to \mathcal{F}_ϕ (see Figure 2). As shown in Table 2, utilizing just the old latent space embedding is useful but no more valuable than using internal adaptation.

In addition to ablations to our own architecture, we also compare our IA component, which can be regarded as a full-rank adaptation scheme, to the low-rank adaptation (LoRA) scheme [Hu et al. 2021]. LoRA typically works well for adapting large language models with a low rank ≤ 8 . However, we did not find any evident improvement over just using LSA when an intrinsic rank of 8 was employed. Even after increasing the rank to 64, the performance gap between the full-rank adaptation scheme and LoRA still remains as listed in Table 2. Though using a low-rank decomposition can reduce the total number of parameters, it increases the computation cost since one more matrix multiplication is needed for each adaptor. Given the small size of our policy network, from our findings we conclude that the full-rank adaptation offers desirable benefits over LoRA.

10 LATENT SPACE ANALYSIS

In this section, we provide more insights on the ability of AdaptNet to successfully control and modify the latent space.

10.1 Latent Space Visualization

Figure 17 visualizes the latent space for different motion style transfer tasks. For each task, a controller was trained using AdaptNet starting from the same pre-trained locomotion policy of walking. During adaptation training here, we use only the latent space injection component as in \mathcal{Z}^0 for all models. We also remove the regularization term \mathcal{I}_ϕ in Eq. 12 and prolong the training time to let AdaptNet fit the style motions as much as possible. After training, we collect samples for each stylized motion from the simulated character following a straight path without any goal-direction changes.

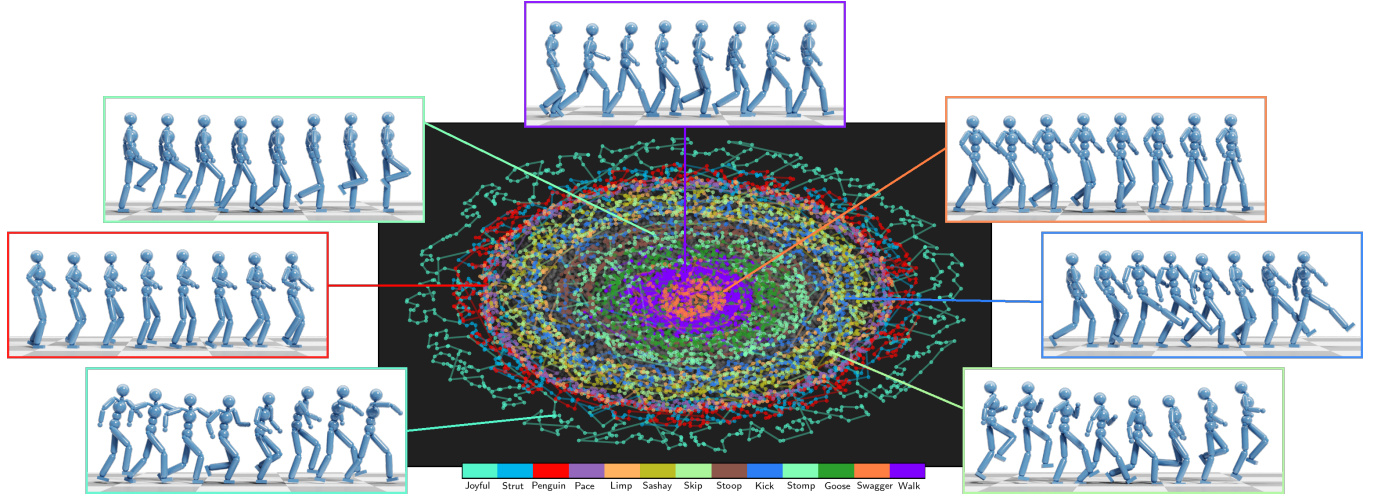


Fig. 17. Latent space visualization with respect to different styles of walk-related motions. The latent representations of the stylized motions are obtained by AdaptNet without using the internal adaptation component. The walk motion (dark purple) near the center is provided by the pre-trained policy based on which AdaptNet performs adaptation to learn the stylized motions. The visualization is achieved using multidimensional scaling technique to project the latent representations from 260 dimensions to 2 dimensions.

We use a multidimensional scaling technique to reduce the dimension of the collected latent samples.

As seen in the figure, the 2D projection of the latent space exhibits a circular shape with the pre-trained walking policy (dark purple) located near the center. There is a clear and roughly continuous transition when the motion style changes from one to the other, which demonstrates the well structured nature of the latent space with the different motion styles. The distribution of the stylized motion in the visualized space is roughly consistent with the imitation error distribution listed in Table 2 when no internal adaptor is employed. Motions with smaller imitation errors are distributed generally closer to the pre-trained policy while *Joyful Walk* (light green) has the largest error and is located the farthest away from the center of the circle. We also note the *Penguin Walk* (red) and *Pace* (light purple) show greater differences in frequency and speed and appear farther away from the center of the figure. This indicates that the distribution in the latent space not only reflects the pose similarity between motions but also some semantic information, like motion rhythm and gait frequency. Similar conclusions have been drawn by recent work in the field of image generation, where the latent space for image generation is considered to capture semantic information more than just simple color transformations [Epstein et al. 2022; Jahanian et al. 2020; Shen et al. 2020].

10.2 Latent Injection Regularization

In Figure 18, we show the latent visualizations of several motions generated by AdaptNet when L2 regularization is applied on the injected latent. For comparison, we highlight in white each motion's distributions in the full latent space shown in Figure 17. In the lower figures, the dark purple points represent the latent embedding of the pre-trained walking, while the gray points are generated by the pre-trained encoder \mathcal{E}_ξ when the simulated character performs stylized motions. Other colors represent varying levels of regularization, as

shown. The goal of regularization is to ensure that the generated latent can fall into the manifold composed of the gray dots. This represents a relatively safe region where the latent space is expected to be handled properly by the pre-trained policy.

In the *Stoop* task, there is almost no difference with and without using the L2 regularization. All visualized samples are overlapped together and covered by the gray region. This is expected given that the style motion of *Stoop* is close to the walking motion in the latent space. However, in the example of *Pace*, there is a clear separation when different regularization coefficients are employed. Note when a coefficient of 0.1 is taken, the generated stylized motion (orange) is overlapped with the walking motion (dark purple). AdaptNet, in this case, is over-regularized. It yields to the pre-trained policy and fails to adapt the pre-trained policy to perform the desired stylized motion. In contrast, without regularization ($\beta = 0$), the latent is already outside of the safe, gray region. AdaptNet, in this case, simply overfits to imitating the style motion and loses the ability to perform goal-steering navigation. While in *Jaunty Skip*, any β -value can be employed, in *Limp* a β -value of 0.01 best ensures that the latent space stays into the grey manifold while attaining high imitation performance. In all adaptation tasks detailed in the paper, we found $\beta = 0.01$ to be sufficient. We note that such regularization is not necessary in other tested adaptation tasks without motion style transfer. In such cases, the new expected motions are close to the original policy and already lie in the safe region. We refer to the supplementary video for a visual comparison of the generated motions when different regularization coefficients are employed.

11 CONCLUSIONS

This paper presents AdaptNet, an approach for adapting existing character RL-based control policies to new motion tasks. Our approach applies two strategies. The first adapts the latent space by conditioning on the character's state and allowing the addition of

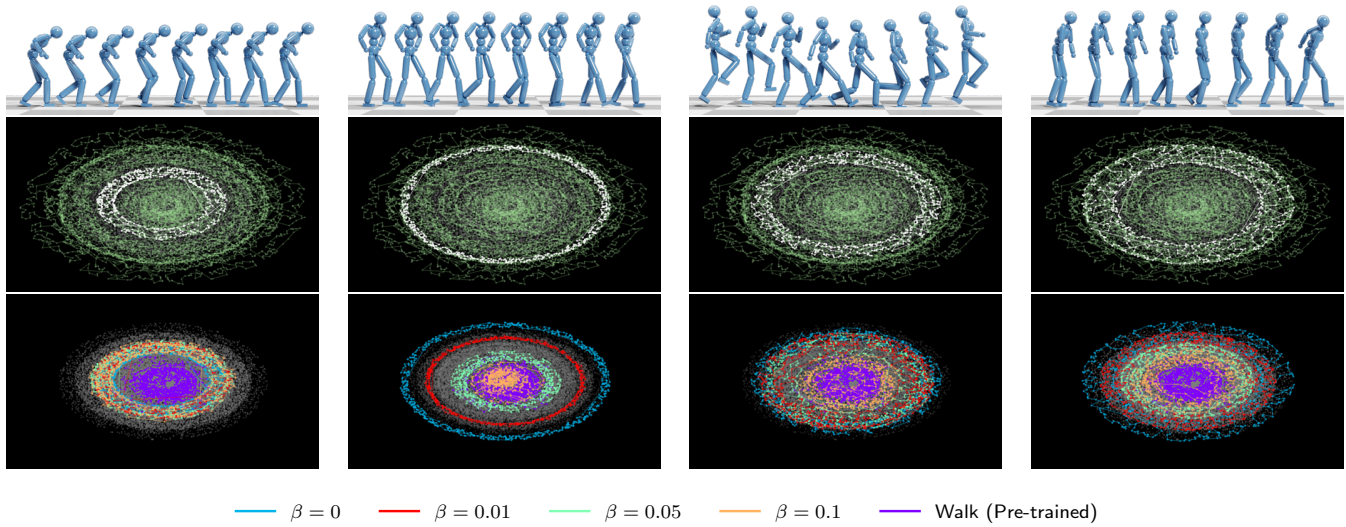


Fig. 18. Latent space distributions of *Stoop*, *Pace*, *Jaunty Walk*, *Limp* (left to right). The top figures show the distribution of the stylized motions in the full latent space without regularization, and the bottom figures show the distribution with regularization applied during adaptation training. Gray points shown in the bottom figures are the latent embeddings generated by the pre-trained encoder \mathcal{E}_ξ while the character performs the stylized motions. β is the regularization coefficient from Eq. 12.

new control inputs that will allow the control policy to perform new tasks. The second aims at control refinement which allows policy adaptation by shifting the original policy and generating new control actions based on new training. Importantly, AdaptNet training always begins with having no (zero) influence, starting from the existing policy and increasing its influence as training proceeds.

We demonstrate that a previously trained control policy for locomotion can be adapted to support diverse style transfer, morphological changes including limb length variation and locked joints, and terrain adaptation including varied friction and geometry. These adaptations are also very efficient to learn. While the original locomotion policy requires 26 hours of training, our style adaptations take less than thirty minutes to produce a full controller that is capable of goal-directed steering while adhering to a specified walking style. More extreme adaptations require more time, but training is still far more efficient than the cost of learning the initial policy.

A core limitation of this work is that policy adaptation requires an existing pre-trained policy, and thus it cannot act to produce new motions on its own. While it is capable of migrating the policy to many new behaviors and conditions, extreme adaptations (e.g., training a jumping action with long flight phase from a walking controller) do not produce the expected results. We believe this is due to the distinct characteristics of the two behaviors and we see such ‘deep’ adaptation as a direction for future work. Also, while we demonstrate smooth interpolation between latent space embeddings when we employ control-layer refinement, interpolation does not always produce coherent in-between behaviors. As we show in Section 9.2, an improper choice of the target latent space could lead to undesired control results. As such, we found starting with a proper latent space is important for obtaining high-quality controllers.

In the current work, we use the recent approach of Xu et al. [2023] for pre-training an initial policy that is then modified by AdaptNet.

In the future, we would like to see how well other recent approaches for training physics-based controllers [Peng et al. 2022, 2021; Yao et al. 2022] can work with our proposed approach. We would also like to investigate how our approach can be extended to generate a well-represented latent space that can be further exploited for motion synthesis. This opens up many avenues for further research, including latent space disentanglement, inversion, and shaping.

ACKNOWLEDGMENTS

We acknowledge the support of the Natural Sciences and Engineering Research Council of Canada (NSERC) and the National Science Foundation under Grants No. IIS-2047632 and IIS-2232066. Support for the first author was made through a generous gift from Roblox. The Bellairs Workshop on Computer Animation was instrumental in the conception of the research presented in this paper.

REFERENCES

- Rameen Abdal, Yipeng Qin, and Peter Wonka. 2019. Image2StyleGAN: How to Embed Images Into the StyleGAN Latent Space?. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4432–4441.
- Kfir Aberman, Yijia Weng, Dani Lischinski, Daniel Cohen-Or, and Baoquan Chen. 2020. Unpaired Motion Style Transfer from Video to Animation. *ACM Transactions on Graphics* 39, 4 (2020).
- Armen Aghajanyan, Sonal Gupta, and Luke Zettlemoyer. 2021. Intrinsic Dimensionality Explains the Effectiveness of Language Model Fine-Tuning. In *59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 7319–7328.
- Ferran Alet, Tomas Lozano-Perez, and Leslie P. Kaelbling. 2018. Modular meta-learning. In *Conference on Robot Learning (Proceedings of Machine Learning Research, Vol. 87)*. 856–868.
- Marcin Andrychowicz, Misha Denil, Sergio Gómez Colmenarejo, Matthew W. Hoffman, David Pfau, Tom Schaul, Brendan Shillingford, and Nando de Freitas. 2016. Learning to Learn by Gradient Descent by Gradient Descent. In *Neural Information Processing Systems*. 3988–3996.
- Kevin Bergamin, Simon Clavet, Daniel Holden, and James Richard Forbes. 2019. DReCon: Data-Driven Responsive Control of Physics-Based Characters. *ACM Transactions on Graphics* 38, 6 (2019).

- David Berthelot, Thomas Schumm, and Luke Metz. 2017. BEGAN: Boundary Equilibrium Generative Adversarial Networks. arXiv:1703.10717 [cs.LG]
- Piotr Bojanowski, Armand Joulin, David Lopez-Pas, and Arthur Szlam. 2018. Optimizing the Latent Space of Generative Networks. In *International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*. 600–609.
- Jason Chemin and Jehee Lee. 2018. A Physics-Based Juggling Simulation Using Reinforcement Learning. In *ACM SIGGRAPH Conference on Motion, Interaction and Games*. Article 3.
- Junyoung Chung, Caglar Gulcehre, Kyunghyun Cho, and Yoshua Bengio. 2014. Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling. In *NIPS 2014 Workshop on Deep Learning*.
- Coline Devin, Abhishek Gupta, Trevor Darrell, Pieter Abbeel, and Sergey Levine. 2017. Learning Modular Neural Network Policies for Multi-Task and Multi-Robot Transfer. In *IEEE International Conference on Robotics and Automation*. 2169–2176.
- Yan Duan, John Schulman, Xi Chen, Peter L. Bartlett, Ilya Sutskever, and Pieter Abbeel. 2016. RL²: Fast Reinforcement Learning via Slow Reinforcement Learning. arXiv:1611.02779 [cs.AI]
- Dave Epstein, Taesung Park, Richard Zhang, Eli Shechtman, and Alexei A. Efros. 2022. BlobGAN: Spatially Disentangled Scene Representations. In *Computer Vision – ECCV 2022*. 616–635.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks. In *International Conference on Machine Learning*. 1126–1135.
- A. Frezzato, A. Tangri, and S. Andrews. 2022. Synthesizing Get-Up Motions for Physics-based Characters. *Computer Graphics Forum* 41, 8 (2022), 207–218.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. 2016. Domain-Adversarial Training of Neural Networks. *Journal of Machine Learning Research* 17, 1 (2016), 2096–2030.
- Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. 2017. Improved Training of Wasserstein GANs. In *Neural Information Processing Systems*, Vol. 30. 5769–5779.
- Abhishek Gupta, Russell Mendonca, YuXuan Liu, Pieter Abbeel, and Sergey Levine. 2018. Meta-Reinforcement Learning of Structured Exploration Strategies. In *Neural Information Processing Systems*, Vol. 31. 5307–5316.
- Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. 2017. Reinforcement Learning with Deep Energy-Based Policies. In *International Conference on Machine Learning*. 1352–1361.
- Tatsuya Harada, Sou Taoka, Taketoshi Mori, and Tomomasa Sato. 2004. Quantitative Evaluation Method for Pose and Motion Similarity Based on Human Perception. In *IEEE/RAS International Conference on Humanoid Robots*, Vol. 1. 494–512.
- Félix G. Harvey, Mike Yurick, Derek Nowrouzezahrai, and Christopher Pal. 2020. Robust Motion In-betweening. *ACM Transactions on Graphics* 39, 4, Article 60 (2020).
- Nicolas Heess, Jonathan J Hunt, Timothy P Lillicrap, and David Silver. 2015. Memory-based control with recurrent neural networks. arXiv:1512.04455 [cs.LG]
- Nicolas Heess, Dhruva TB, Srinivasan Sriram, Jay Lemmon, Josh Merel, Greg Wayne, Yuval Tassa, Tom Erez, Ziyu Wang, S. M. Ali Eslami, Martin Riedmiller, and David Silver. 2017. Emergence of Locomotion Behaviours in Rich Environments. arXiv:1707.02286 [cs.AI]
- Donald Hejira, Lerrel Pinto, and Pieter Abbeel. 2020. Hierarchically Decoupled Imitation For Morphological Transfer. In *37th International Conference on Machine Learning*, Vol. 119. 4159–4171.
- Jonathan Ho and Stefano Ermon. 2016. Generative Adversarial Imitation Learning. *Advances in Neural Information Processing Systems* 29 (2016).
- Rein Houthoofd, Yuhua Chen, Phillip Isola, Bradley Stadie, Filip Wolski, OpenAI Jonathan Ho, and Pieter Abbeel. 2018. Evolved policy gradients. In *Neural Information Processing Systems*, Vol. 31. 5405–5414.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. LoRA: Low-Rank Adaptation of Large Language Models. arXiv:2106.09685 [cs.CL]
- Ali Jahanian, Lucy Chai, and Phillip Isola. 2020. On the “Steerability” of Generative Adversarial Networks. In *International Conference on Learning Representations*.
- Jordan Juravsky, Yunrong Guo, Sanja Fidler, and Xue Bin Peng. 2022. PADL: Language-Directed Physics-Based Character Control. In *SIGGRAPH Asia 2022 Conference Papers*. Article 19.
- Andrej Karpathy and Michiel van de Panne. 2012. Curriculum Learning for Motor Skills. In *Canadian Conference on Artificial Intelligence*. Springer, 325–330.
- Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. 2020. Analyzing and Improving the Image Quality of StyleGAN. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8110–8119.
- Diederik P. Kingma and Jimmy Ba. 2017. Adam: A Method for Stochastic Optimization. arXiv:1412.6980 [cs.LG]
- Ariel Kwiatkowski, Eduardo Alvarado, Vicky Kalogeiton, C. Karen Liu, Julien Pettré, Michiel van de Panne, and Marie-Paule Cani. 2022. A Survey on Reinforcement Learning Methods in Character Animation. *Computer Graphics Forum* 41, 2 (2022), 613–639.
- Chunyan Li, Heerad Farkhoor, Rosanne Liu, and Jason Yosinski. 2018. Measuring the Intrinsic Dimension of Objective Landscapes. In *International Conference on Learning Representations*.
- Jae Hyun Lim and Jong Chul Ye. 2017. Geometric GAN. arXiv:1705.02894 [stat.ML]
- Hung Yu Ling, Fabio Zinno, George Cheng, and Michiel van de Panne. 2020. Character controllers using motion VAEs. *ACM Transactions on Graphics* 39, 4, Article 40 (2020).
- Libin Liu and Jessica Hodgins. 2017. Learning to Schedule Control Fragments for Physics-Based Characters Using Deep Q-Learning. *ACM Transactions on Graphics* 36, 4, Article 42a (2017).
- Libin Liu and Jessica Hodgins. 2018. Learning Basketball Dribbling Skills Using Trajectory Optimization and Deep Reinforcement Learning. *ACM Transactions on Graphics* 37, 4, Article 142 (2018), 14 pages.
- Yunhao Luo, Kaixiang Xie, Sheldon Andrews, and Paul Kry. 2021. Catching and Throwing Control of a Physically Simulated Hand. In *ACM SIGGRAPH Conference on Motion, Interaction and Games*. Article 15.
- Viktor Makovychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, and Gavriel State. 2021. Isaac Gym: High Performance GPU-Based Physics Simulation For Robot Learning. arXiv:2108.10470 [cs.RO]
- I. Mason, S. Starke, H. Zhang, H. Bilen, and T. Komura. 2018. Few-shot Learning of Homogeneous Human Locomotion Styles. *Computer Graphics Forum* 37, 7 (2018), 143–153.
- Josh Merel, Yuval Tassa, Dhruva TB, Sriram Srinivasan, Jay Lemmon, Ziyu Wang, Greg Wayne, and Nicolas Heess. 2017. Learning human behaviors from motion capture by adversarial imitation. arXiv:1707.02201 [cs.RO]
- Josh Merel, Saran Tunyasuvunakool, Arun Ahuja, Yuval Tassa, Leonard Hasenclever, Vu Pham, Tom Erez, Greg Wayne, and Nicolas Heess. 2020. Catch & Carry: Reusable Neural Controllers for Vision-Guided Whole-Body Tasks. *ACM Transactions on Graphics* 39, 4, Article 39 (2020).
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. 2013. Playing Atari with Deep Reinforcement Learning. arXiv:1312.5602 [cs.LG]
- Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaohu Qie. 2023. T2I-Adapter: Learning Adapters to Dig out More Controllable Ability for Text-to-Image Diffusion Models. arXiv:2302.08453 [cs.CV]
- Alex Nichol, Joshua Achiam, and John Schulman. 2018. On First-Order Meta-Learning Algorithms. arXiv:1803.02999 [cs.LG]
- Emilio Parisotto, Lei Jimmy Ba, and Ruslan Salakhutdinov. 2016. Actor-Mimic: Deep Multitask and Transfer Reinforcement Learning. In *International Conference on Learning Representations*.
- Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. 2018a. Deep-Mimic: Example-Guided Deep Reinforcement Learning of Physics-Based Character Skills. *ACM Transactions on Graphics* 37, 4, Article 143 (2018).
- Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. 2018b. Sim-to-Real Transfer of Robotic Control with Dynamics Randomization. In *IEEE International Conference on Robotics and Automation*. 3803–3810.
- Xue Bin Peng, Michael Chang, Grace Zhang, Pieter Abbeel, and Sergey Levine. 2019. MCP: Learning Composable Hierarchical Control with Multiplicative Compositional Policies. In *Advances in Neural Information Processing Systems*. 3681–3692.
- Xue Bin Peng, Yunrong Guo, Lina Halper, Sergey Levine, and Sanja Fidler. 2022. ASE: Large-Scale Reusable Adversarial Skill Embeddings for Physically Simulated Characters. *ACM Transactions on Graphics* 41, 4, Article 94 (2022).
- Xue Bin Peng, Ze Ma, Pieter Abbeel, Sergey Levine, and Angjoo Kanazawa. 2021. AMP: Adversarial Motion Priors for Stylized Physics-Based Character Control. *ACM Transactions on Graphics* 40, 4, Article 144 (2021).
- Phil Pope, Chen Zhu, Ahmed Abdelkader, Micah Goldblum, and Tom Goldstein. 2021. The Intrinsic Dimension of Images and Its Impact on Learning. In *International Conference on Learning Representations*.
- Alec Radford, Luke Metz, and Soumith Chintala. 2016. Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks. arXiv:1511.06434 [cs.LG]
- Aravind Rajeswaran, Sarveer Ghotra, Balaraman Ravindran, and Sergey Levine. 2017. EPOpt: Learning Robust Neural Network Policies Using Model Ensembles. In *International Conference on Learning Representations*.
- Sachin Ravi and Hugo Larochelle. 2017. Optimization as a Model for Few-Shot Learning. In *International Conference on Learning Representations*.
- Andrei A. Rusu, Sergio Gomez Colmenarejo, Caglar Gulcehre, Guillaume Desjardins, James Kirkpatrick, Razvan Pascanu, Volodymyr Mnih, Koray Kavukcuoglu, and Raia Hadsell. 2016a. Policy Distillation. arXiv:1511.06295 [cs.LG]
- Andrei A. Rusu, Neil C. Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. 2016b. Progressive Neural Networks. arXiv:1606.04671 [cs.LG]
- Andrei A Rusu, Matej Večerík, Thomas Rothörl, Nicolas Heess, Razvan Pascanu, and Raia Hadsell. 2017. Sim-to-Real Robot Learning from Pixels with Progressive Nets. In *Conference on Robot Learning*. 262–270.

- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal Policy Optimization Algorithms. arXiv:1707.06347 [cs.LG]
- Yujun Shen, Jinjin Gu, Xiaou Tang, and Bolei Zhou. 2020. Interpreting the Latent Space of GANs for Semantic Face Editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9243–9252.
- Tom Silver, Kelsey Allen, Josh Tenenbaum, and Leslie Kaelbling. 2019. Residual Policy Learning. arXiv:1812.06298 [cs.RO]
- Sebastian Starke, Ian Mason, and Taku Komura. 2022. DeepPhase: Periodic Autoencoders for Learning Motion Phase Manifolds. *ACM Transactions on Graphics* 41, 4, Article 136 (2022).
- Jeff KT Tang, Howard Leung, Taku Komura, and Hubert PH Shum. 2008. Emulating human perception of motion similarity. *Computer Animation and Virtual Worlds* 19, 3–4 (2008), 211–221.
- Tianxin Tao, Matthew Wilson, Ruiyu Gou, and Michiel van de Panne. 2022. Learning to Get Up. In *ACM SIGGRAPH 2022 Conference Proceedings*. Article 47.
- Chen Tessler, Yoni Kasten, Yunrong Guo, Shie Mannor, Gal Chechik, and Xue Bin Peng. 2023. CALM: Conditional Adversarial Latent Models for Directable Virtual Characters. In *ACM SIGGRAPH 2023 Conference Proceedings*. Article 37.
- Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. 2017. Adversarial Discriminative Domain Adaptation. In *IEEE Conference on Computer Vision and Pattern Recognition*. 2962–2971.
- Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno A. Olshausen, and Trevor Darrell. 2021. Tent: Fully Test-Time Adaptation by Entropy Minimization. In *International Conference on Learning Representations*.
- Jungdam Won, Deepak Gopinath, and Jessica Hodgins. 2021. Control Strategies for Physically Simulated Characters Performing Two-Player Competitive Sports. *ACM Transactions on Graphics* 40, 4, Article 146 (2021).
- Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. 2016. Learning a Probabilistic Latent Space of Object Shapes via 3D Generative-Adversarial Modeling. In *Advances in Neural Information Processing Systems*, Vol. 29.
- Zhaoming Xie, Hung Yu Ling, Nam Hee Kim, and Michiel van de Panne. 2020. ALL-STEPS: Curriculum-driven Learning of Stepping Stone Skills. *Computer Graphics Forum* 39, 8 (2020), 213–224.
- Zhaoming Xie, Sebastian Starke, Hung Yu Ling, and Michiel van de Panne. 2022. Learning Soccer Juggling Skills with Layer-Wise Mixture-of-Experts. In *ACM SIGGRAPH 2022 Conference Proceedings*. Article 25.
- Pei Xu and Ioannis Karamouzas. 2021. A GAN-Like Approach for Physics-Based Imitation Learning and Interactive Character Control. *Proceedings of the ACM on Computer Graphics and Interactive Techniques* 4, 3, Article 44 (2021).
- Pei Xu, Xiumin Shang, Victor Zordan, and Ioannis Karamouzas. 2023. Composite Motion Learning with Task Control. *ACM Transactions on Graphics* 42, 4, Article 93 (2023).
- Zhongwen Xu, Hado P van Hasselt, and David Silver. 2018. Meta-Gradient Reinforcement Learning. In *Advances in Neural Information Processing Systems*, Vol. 31.
- Heyuan Yao, Zhenhua Song, Baoquan Chen, and Libin Liu. 2022. ControlVAE: Model-Based Learning of Generative Controllers for Physics-Based Characters. *ACM Transactions on Graphics* 41, 6, Article 183 (2022).
- Zhiqi Yin, Zeshi Yang, Michiel van de Panne, and Kangkang Yin. 2021. Discovering Diverse Athletic Jumping Strategies. *ACM Transactions on Graphics* 40, 4, Article 91 (2021).
- Wenhao Yu, Greg Turk, and C Karen Liu. 2018. Learning Symmetric and Low-Energy Locomotion. *ACM Transactions on Graphics* 37, 4, Article 144 (2018).
- Lvmin Zhang and Maneesh Agrawala. 2023. Adding Conditional Control to Text-to-Image Diffusion Models. arXiv:2302.05543 [cs.CV]
- Peiye Zhuang, Oluwasanmi O Koyejo, and Alex Schwing. 2021. Enjoy Your Editing: Controllable GANs for Image Editing via Latent Space Navigation. In *International Conference on Learning Representations*.

A LOCOMOTION TASK SETUP

The goal state $\mathbf{g}_t \in \mathbb{R}^4$ includes a 2D unit vector representing the direction to the target location, the horizontal distance from the character to the goal, i.e., $\|\mathbf{x}_t^{\text{root}} - \mathbf{p}_{\text{goal}}\|$, and the preferred speed $\|\mathbf{v}_t^*\|$. The preferred speed is sampled from $[1, 1.5]$ in the unit of m/s for crouching and walking motions, and from $[1, 3]$ for running. The goal direction is sampled from $[0, 2\pi)$. A timer variable is sampled from $[3, 5]$ in the unit of s for walking motions, and from $[2, 3]$ for running. We use these three goal variables to obtain the target location. As such, we can perform speed control during the location

Table 3. Training hyperparameters.

Parameter	Value
policy network learning rate	5×10^{-6}
critic network learning rate	1×10^{-4}
discriminator learning rate	1×10^{-5}
reward discount factor (γ)	0.95
GAE discount factor (λ)	0.95
surrogate clip range (ϵ)	0.2
gradient penalty coefficient (λ^{GP})	10
number of PPO workers (simulation instances)	512
PPO replay buffer size	4096
PPO batch size	256
PPO optimization epochs	5
discriminator replay buffer size	8192
discriminator batch size	512
Regularizer for internal adaptors (β)	0.01
Regularizer coefficient for latent space adaptor (κ)	0.01

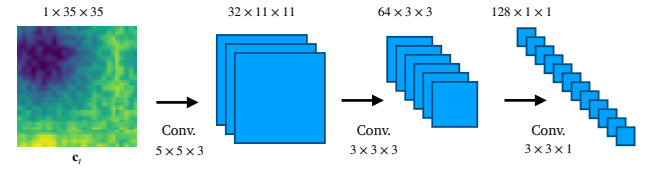


Fig. 19. Network structure of the heightmap encoder \mathcal{G}_ϕ giving additional control input \mathbf{c}_t . The convolution operation Conv. has format of kernel height \times kernel width \times stride. The dimension of \mathbf{c}_t and the feature maps are shown in the format of channels \times height \times width.

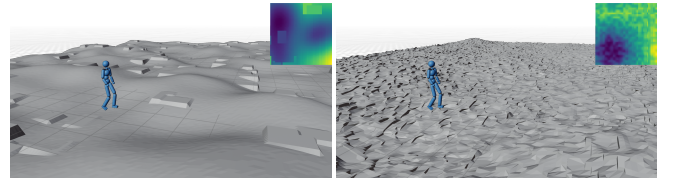


Fig. 20. Procedurally generated terrains with minimaps showing the local heightmaps used as additional control input \mathbf{c}_t . The 35×35 heightmaps encode a $3.4 \text{ m} \times 3.4 \text{ m}$ area centered around the character.

targeting. The goal-directed task reward is

$$r_t = \begin{cases} \exp(-3\|\dot{\mathbf{x}}_{t+1}^{\text{root}}/T - \mathbf{v}_t^*\|^2 / \|\mathbf{v}_t^*\|^2) & \text{if } \|\mathbf{x}_{t+1} - \mathbf{p}_{\text{goal}}\| > R \\ 1 & \text{otherwise,} \end{cases}$$

where $R = 0.5$ is the goal radius of the target location, $T = 1/30 \text{ s}$ is the time interval between two frames, $\dot{\mathbf{x}}_{t+1}^{\text{root}}/T$ denotes the horizontal velocity of the root link at time $t + 1$, and \mathbf{v}_t^* is the target velocity toward the goal location.

B HYPERPARAMETERS

The hyperparameters used for policy training are listed in Table 3. Half of the training samples used by the discriminator are drawn from the simulated character and half are sampled from the reference motions. The learning objective weight ω_k in Equation 9 of the main text is 0.5 for imitation and 0.5 for the task of goal-steering

navigation. During adaptation, we use 0.35 for imitation and 0.65 for navigation in motion style transfer tasks and 0.5 for both imitation and navigation in other tasks.

C TERRAIN MAPS

Terrains are generated using a combination of Perlin noise and procedural generation. To generate more challenging terrains, we add extra uniform noise to make the terrain more uneven, while including some randomly generated rectangular blocks for more stepped terrain. In the most rugged example, we generate terrain by

applying uniform noise on a terrain composed of multiple surfaces with slopes varying from -0.3 to 0.3. An extra slope threshold is applied to cut off the bulges on the terrain and make it rugged. The generated terrains in all examples are normalized to have a maximal height of 0.75 m and a minimum height of -0.75 m. Figure 19 shows the network structure of the heightmap encoding module \mathcal{G}_ϕ , which has three convolutional layers. Following previous work [Mnih et al. 2013], we do not use any pooling layer in the network. Resulting terrain examples appear in Figure 20.