# Learning to Brachiate via Simplified Model Imitation

Daniele Reda*
dreda@cs.ubc.ca
University of British Columbia
Vancouver, Canada

Hung Yu Ling*
hyuling@cs.ubc.ca
University of British Columbia
Vancouver, Canada

Michiel van de Panne
van@cs.ubc.ca
University of British Columbia
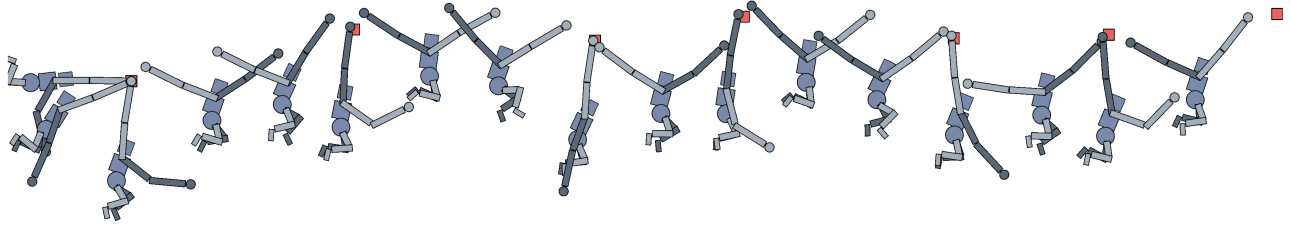Vancouver, Canada

**Figure 1: We use a two-stage simplified model imitation learning approach to produce realistic physics-based brachiation motions. The learned control policies can traverse challenging handhold sequences and can demonstrate emergent *pumping* behavior to build momentum for large gaps.**

## ABSTRACT

Brachiation is the primary form of locomotion for gibbons and siamangs, in which these primates swing from tree limb to tree limb using only their arms. It is challenging to control because of the limited control authority, the required advance planning, and the precision of the required grasps. We present a novel approach to this problem using reinforcement learning, and as demonstrated on a finger-less 14-link planar model that learns to brachiate across challenging handhold sequences. Key to our method is the use of a simplified model, a point mass with a virtual arm, for which we first learn a policy that can brachiate across handhold sequences with a prescribed order. This facilitates the learning of the policy for the full model, for which it provides guidance by providing an overall center-of-mass trajectory to imitate, as well as for the timing of the holds. Lastly, the simplified model can also readily be used for planning suitable sequences of handholds in a given environment. Our results demonstrate brachiation motions with a variety of durations for the flight and hold phases, as well as emergent extra back-and-forth swings when this proves useful. The system is evaluated with a variety of ablations. The method enables future work towards more general 3D brachiation, as well as using simplified model imitation in other settings. For videos, supplementary material and code, visit: https://brachiation-rl.github.io/brachiation.

## CCS CONCEPTS

• **Computing methodologies** → *Physical simulation*; *Reinforcement learning*.

## KEYWORDS

physics-based character animation, motion planning, reinforcement learning

## 1 INTRODUCTION

Brachiation is the form of locomotion that consists of moving between tree branches or handholds using the arms alone [Fleagle 1974; Preuschoft and Demes 1985]. Gibbons and siamangs make this seem effortless and are among the world's most agile brachiators. Brachiating movements can be classified into two types: *continuous contact brachiation*, characterized by the animal mantaining contact at all time with a handhold such as a tree branch, and *ricochetal brachiation*, which involves a flight phase between successive grasp [Bertram et al. 1999; Wan et al. 2015]. These two modes of brachiation are analogous to walking and running, respectively, with ricochetal brachiation used at faster speeds.

Machine learning techniques such as deep reinforcement learning (RL) have been previously applied to legged, aerial, and underwater locomotion [Luo et al. 2020; Min et al. 2019; Won et al. 2018; Xie et al. 2020]. Despite the apparent similarity to these types of locomotion, brachiation is unique in a number of ways. First, the choice of handholds is often discrete, i.e., the choice of particular branches, and therefore unlike the continuous terrain usually available for legged locomotion. This offers apparently fewer features to control, i.e., less control authority, while at the same time necessitating high spatial precision to effect a grasp. Second, advance

planning of the motion across a sequence of handholds becomes more important because of the momentum that may be needed to reach future handholds. Lastly, brachiation offers the possibility of very efficient horizontal motion; an ideal point mass can follow alternating pendulum-like swings and parabolic flight phases with no loss of energy [Bertram et al. 1999].

In our work, we present a fully learned solution for simulating a 14-link planar articulated model. Similar to previous work on brachiation, our model uses a pinned attachment model as a proxy for the grasping mechanics. The learned control policy demonstrates ricochetal brachiation across challenging sequences of handholds that, to the best of our knowledge, is the most capable of its kind.

Our primary contributions are as follows:

- We present an effective learning-based solution for planar brachiation. The learned control policies can traverse challenging sequences of handholds using ricochetal motions with significant flight phases, and demonstrate emergent phenomena, such as extra back-and-forth swings when needed.

- We propose a general two-stage RL method for leveraging simplified physical models. RL on the simplified model allows for fast and efficient exploration, quickly producing an effective control policy. The full model then uses imitation learning of the motions resulting from the simplified-model for efficient learning of complex behaviors.

## 2  RELATED WORK

Our work sits at the intersection of physics-based character animation, the study of brachiation, and a broad set of methods that leverage abstract models in support of motion planning and control.

### 2.1  Brachiation

Due to the fascinating and seemingly effortless nature of brachiation, there exists a long history of attempts to reproduce this often-graceful type of motion, both in simulation and on physical robots. The key challenge lies in developing control strategies that are well suited to the underactuated and dynamic nature of the task.

A starting point is to first study brachiation in nature, such as via a detailed film study [Fleagle 1974]. Among other insights, they note that the siamang is able to limit lateral motion of the center of mass between handholds, in part due to extensive rotation at the wrist, elbow, and shoulder. A benefit of longer arms for reduced energetic costs of brachiation can also be found [Preuschoft and Demes 1985]. A more recent study [Michilsens et al. 2009] notes that "compared with other primates, the elbow flexors of gibbons are particularly powerful, suggesting that these muscles are particularly important for a brachiating lifestyle." We note that the limited lateral motion and strong elbow flexors correspond well to the capabilities afforded by sagittal-plane point-mass models.

One of the earliest results for brachiation capable of varying height and distance uses an intricate heuristic procedure and manually structured motion phases, as applied to a two-link robot with no flight phases [Saito 1992; Saito et al. 1994]. A loosely comparable approach is proposed for a 3-link planar model, restricted to horizontal motions [Zhang and Wong 1999]. A point mass model

can be shown to produce solutions with zero energetic cost for both continuous contact and ricochetal brachiation, for horizontal brachiation with regular spacing of handholds [Bertram et al. 1999]. The solutions are based on circular pendular motion, alternating with parabolic free flight phases in the case of ricochetal motion. It is found that natural gibbon motion is even smoother than the motions predicted by this model, particularly for handhold forces. Several results successfully developed horizontal continuous brachiation strategies for a two-link robot [Nakanishi et al. 2000] and motions for a three-link model as demonstrated on a humanoid robot [Kajima et al. 2004]. Zero-energy costs are later demonstrated for a five-link model [Gomes and Ruina 2005].

Methods and results since 2015 have continued to make additional progress. Non-horizontal ricochetal brachiation is demonstrated for a 2-link primate robot [Wan et al. 2015] and uses an analytic solver for the boundary conditions of each swing and a Lyapunov-based tracking controller. A method is developed for planar brachiation on flexible cables [Farzan et al. 2018], for a 3-link model and corresponding robot. A multiple shooting method is used to plan open-loop feedforward commands. Brachiation control of a 3-link robot is also designed and demonstrated using iLQR [Yang et al. 2019], for animated sequences of single swings. A pendulum model and a planar articulated model are described in the PhD work of Berseth [Berseth 2019], based on a finite-state-machine control structure. While demonstrating potential, the results are based on a "grab anywhere" model and produce controllers with limited capabilities and realism. The control for a transverse (sideways) ricochetal brachiation has also been investigated recently, for a 4-link arms, body, and tail model [Lin and Yang 2020].

Our work differs from the bulk of the above work in the following ways. We demonstrate the ability to learn planar ricochetal brachiation with minimal manual engineering, and that can traverse handhold sequences with significant variation. We demonstrate emergent behaviors, such as additional back-and-forth swings as may be needed to proceed in some situations. Our approach leverages the benefits of learning with a simplified model, showing that the simplified-model control policies can provide highly-effective guidance for learning with the full articulated-body model. The learned control policies anticipate upcoming handholds, and can be readily integrated with a simple motion planner.

### 2.2  Physics-based Character Skills

Physics-based worlds require simulated characters that are capable of many skills, including the ability to move through all kinds of environments with speed and agility. Significant progress has been made for over thirty years, although it is only recently that we have begun to see approaches that demonstrate scalability, in particular with regard to being able to imitate a wide variety of motion capture clips with a physics-based human model, e.g., [Bergamin et al. 2019; Chentanez et al. 2018; Merel et al. 2018; Park et al. 2019; Peng et al. 2018; Wang et al. 2020; Won et al. 2020]. Recent work has also demonstrated the ability to learn bipedal locomotion that are capable of navigating difficult step sequences, with learning curricula and in the absence of motion capture data [Xie et al. 2020].

## 2.3 Learning with Simplified Models

Simplified physical approximations are a standard way to create models that are tractable to simulate and control, particularly for legged locomotion and brachiation. A common strategy is to first create a motion plan using the simplified model, e.g., a spring-loaded inverted pendulum [Poulakakis and Grizzle 2009] or centroidal dynamics models, e.g., [Green et al. 2021; Tsounis et al. 2020; Xie et al. 2021], and then treat this as a reference motion to be tracked online [Siciliano et al. 2008]. Both the motion planning and tracking problems have commonly been solved using model-based trajectory optimization (TO) or model-free reinforcement learning (RL) methods. In what follows below, we will use a 'plan+tracking' notation to loosely categorize combinations of motion planning and tracking of the resulting motion plan. There also exists an abundance of literature on pure TO and RL solutions, which we place outside the scope of our discussion, although for evaluation purposes, we do consider an RL-only baseline method.

A prominent example of a TO+TO approach can be found in the control of the Boston Dynamics Atlas robot [Kuindersma 2020], where the planning TO uses a long horizon and solves from scratch, and the tracking TO performs online MPC-style solvers over a short horizon and leverages the planning TO results to warm start the optimization. Another TO+TO example first plans legged locomotion with a sliding-based inverted pendulum, which then initializes a more detailed TO optimization with contact locations and a centroidal model [Kwon et al. 2020], and, lastly, maps this to a full body model using momentum-mapped inverse kinematics. Much recent work in imitation-based physics-based animation and robotics could be categorized as fixed+RL, where captured motion data effectively serves as a fixed motion plan, e.g., [Peng et al. 2018] and related works.

Quadruped control policies have been developed using both TO+RL and RL+TO approaches, e.g., [Gangapurwala et al. 2021] and [Gangapurwala et al. 2020; Xie et al. 2021], respectively. A recent RL+RL approach uses a simplified model (and related policy) that is allowed to sample in state space rather than action space [Sontakke and Ha 2021]. This allows it to behave more like a motion planner and assumes that the result can still provide a useful motion plan to for the full model to imitate. The method demonstrates the ability to avoid a stand-in-place local minima for a simulated quadruped, which may arise due to sparse rewards. Our work develops an RL+RL approach for the challenging control problem of brachiation. In our case, both the simplified and full MDPs work directly in a physically-grounded action space. We demonstrate the benefit of this two-level approach via multiple baselines and ablations.

## 3 ENVIRONMENTS

### 3.1 Simplified Model

Our simplified model treats the gibbon as a point mass equipped with a virtual extensible arm. The virtual arm is capable of grabbing any handhold within the annular area defined by a radius $r \in [r_{min}, r_{max}]$, with $r_{min}$=10 cm, and $r_{max}$=75 cm. We note that while the $r_{max}$ is greater than the arm length of the full model (60 cm), the simplified model is intended to be a proxy for the center of mass of the gibbon, and not the shoulder.

**Table 1: Physical Properties of Full Gibbon Model. The waist joint link mass and length includes the torso and the head. Others include only the direct child link.**

| Joint | Max. Torque (Nm) | Link Mass (kg) | Link Length (cm) |
|---|---|---|---|
| Waist | 21 | 3.69 | 31 |
| Shoulder | 35 | 0.74 | 26 |
| Elbow | 28 | 0.74 | 26 |
| Wrist | 0 | 0.40 | 8 |
| Hip | 14 | 0.37 | 12 |
| Knee | 14 | 0.31 | 12 |
| Ankle | 7 | 0.47 | 8 |

The simplified model dynamics consist of two alternating phases. The swing phase starts when the character performs a grab when one of its hands is at (or nearby) to the target handhold. During this phase, the character can apply a force along the direction of the grabbing arm, either to pull towards or push away from the current handhold. The swing phase dynamics is equivalent to a spring-and-damper point-mass pendulum system. The character can influence the angular velocity by shortening or lengthening the swing arm of the pendulum. The flight phase is defined as the period when neither hand is grabbing. During this phase, the character's trajectory is governed by passive physics, following a parabolic trajectory. The control authority that remains during a flight phase comes from the decision of when to grab onto the next handhold, if it is within reaching distance. The control parameters are described in detail in Section 3.4.

The minimum and maximum arm length are enforced by applying a force impulse when the length constraint is violated during a swing phase. We implement the physics simulation of the simplified environment in PyTorch [Paszke et al. 2019], which can be trivially parallelized on a GPU.

### 3.2 Full Model

The full model is a planar articulated model of a gibbon-like character that approximates the morphology of real gibbons, as reported in [Michilsens et al. 2009], and is simulated using PyBullet [Coumans and Bai 2022].

Our gibbon model consists of 13 hinge joints including shoulders, elbows, wrists, hips, knees, ankles, and a single waist joint at the pelvis. All joints have associated torques, except for the wrist joints, which we consider to be passive. To capture the 3D motion of the ball-and-socket joints in two dimensions, we model the shoulder joints as hinge joints without any corresponding joint limits. This allows the simulated character to produce the *under-swing* motions of real gibbons. Our gibbon has a total mass of 9.8 kg, an arm-reach of 60 cm, and a standing height of 63 centimeters. We artificially increase the mass of the hands and the feet to improve simulation stability. The physical properties of our gibbon model are summarized in Table 1.

The grab behavior in the full model is simulated using point-to-point constraints. This allows us to abstract away the complexity of modelling hand-object dynamics, yet still provides a reasonable approximation to the underlying physics. In our implementation, a point-to-point constraint is created when the character performs
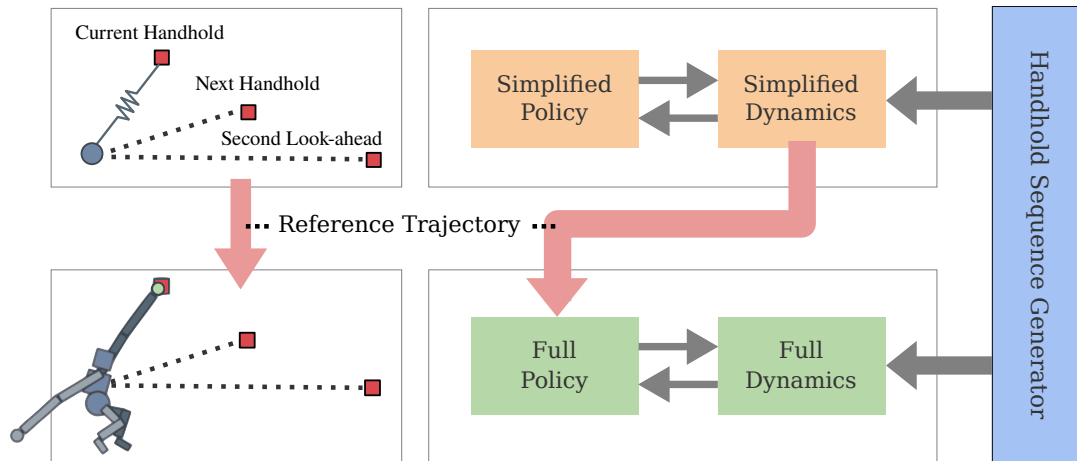
**Figure 2: Overview of our simplified model imitation learning system. The simplified model allows for efficient exploration and quickly produces physically-feasible reference trajectories for the full model to imitate.**

a grab and the hand is within five centimeters of the target handhold. The constraint *pins* the grabbing hand in its current position, as opposed to the target handhold location, to avoid introducing undesirable fictitious forces. At release, the constraint is removed.

### 3.3   Handhold Sequences Generation

Both the simplified and full characters operate in the same environment, which contains handhold sequences. In a brachiation task, the goal is to grab successive handholds precisely and move forward. Successive handholds are generated from uniform distributions defined by the distance, $d \sim U(1, 2)$ meters, and pitch, $\phi \sim [-15°, 15°]$, relative to the previous handhold.

### 3.4   Action Spaces

The simplified and full models both use a control frequency of 60 Hz and a simulation frequency is 480 Hz.

The action space for the simplified model consists of two values, a target length offset and a flag indicating to grab or to release. At the beginning of each control step, the target length offset is added to the current arm length to form the desired target length. In addition, the grab flag is used to determine the current character state. At each simulation step, the applied force is computed using PD control based on the desired target length. We clamp the maximum applied force to 240 N, close to the maximum observed forces in real gibbon reported in [Michilsens et al. 2009].

In the full model, the action space is analogous to that of the simplified model, only extended to every joint. An action consists of 15 parameters representing the joint angle offsets for each of the 13 joints and a flag indicating grab or release for each hand. The applied torques in each joint are computed using PD control. We set the proportional gain, $k_p$, for each joint to be equal to its range of motion in radians and the derivative gain to be $k_p/10$.

### 3.5   State Spaces

Both environments have a similar state space, composed of information of the character and information on the future handholds sequence to grab.

In the simplified environment, the character state is three dimensional consisting of the root velocity in world coordinates and the elapsed time since the start of the current swing phase. The elapsed time information is necessary to make the environment fully observable since the environment enforces a minimum and a maximum grab duration (§4.4). We normalize the elapsed time by dividing by the maximum allowable grab duration. A zero indicates that the character is not grabbing and a one indicates it has grabbed for the maximum amount of time.

The character state for the full model is a 45D vector consisting of root velocity, root pitch, joint angles, joint velocities, grab arm height, and grab states. The torso is used as the root link. The root velocity ($\mathbb{R}^2$) is the velocity in the sagittal plane. Joint angles ($\mathbb{R}^{26}$) are represented as $\cos(\theta)$ and $\sin(\theta)$ of each joint, where $\theta$ is the current angle. Joint velocities ($\mathbb{R}^{13}$) are the angular velocities measured in radians per second. Grab arm height ($\mathbb{R}$) is the distance between the free hand and the root in the upward direction, which can be used to infer whether the next target is reachable. Lastly, the grab states ($\mathbb{R}^2$) are Boolean flags indicating if each hand is currently grabbing.

In both environments, the control policy receives task information pertinent to the handhold locations in addition to the character state. In particular, the task information includes the location of the current handhold and the $N$ upcoming handholds in the character's coordinate frame. We experiment with different values of $N$ in Section 5.1.

## 4   LEARNING CONTROL POLICIES

We use deep reinforcement learning (DRL) to learn brachiation skills in both the simplified and the detailed environment. In RL, at each time step $t$, the control policy reacts to an environment state $s_t$ by performing an action $a_t$. Based on the action performed,

the policy receives a reward signal $r_t = r(s_t, a_t)$ as feedback. In DRL, the policy computes $a_t$ using a neural network $\pi_\theta(a|s)$, where $\pi_\theta(a|\cdot)$ is the probability density of $a$ under the current policy. The goal of DRL is to find the network parameters $\theta$ which maximize the following:

$$J_{RL}(\theta) = E\left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)\right]. \tag{1}$$

Here, $\gamma \in [0, 1)$ is the discount factor so that the sum converges. We solve this optimization problem using the proximal policy optimization (PPO) algorithm [Schulman et al. 2017], a policy gradient actor-critic algorithm. We choose PPO for its effective utilization of hardware resources in parallelized settings which results in reduced wall-clock learning time. Our learning pipeline is based on a publicly available implementation of PPO [Kostrikov 2018]. A review of the PPO algorithm is provided in Appendix B.

## 4.1 System Overview

We start by providing an overview of our system, shown in Figure 2. Our system contains three distinct components: the simplified model, the full model, and a handhold sequence generator. During training, the simplified model is trained first on randomly sampled handhold sequences generated by the handhold generator. After the simplified model is trained, we obtain reference trajectories from the simplified environment on a fixed set of handholds sequences. The full model is subsequently trained by imitating the reference trajectories from the simplified model while optimizing task and style objectives. At evaluation time, the simplified model can be used as a planner to provide guidance for the full model either upfront for an entire handhold sequence or on a per-grab basis. Both models working in tandem allows our system to traverse difficult handholds sequences that would otherwise be challenging for standard reinforcement learning policies.

## 4.2 Learning Simplified Policies

The simplified policy is trained to optimize a sparse reward that is given when the character grabs the next handhold. The fact that the simplified policy can be trained with only the sparse task reward is desirable as it allows for interesting behaviors to emerge. In DRL, reward shaping is often required for the learning to succeed or to significantly speed up the learning. At the same time, reward shaping can also introduce biases that cause the control policy to optimize for task-irrelevant objectives and deviate from the main task. Directly optimizing a policy using the sparse task reward allows us to find solution modes that would otherwise be exceedingly difficult to find.

*Simplified Policy Networks.* The simplified policy contains two neural networks, the controller and the value function. The controller and the value function have similar architecture differing only in the output layer. Each neural network is a three-layered feed-forward network with 256 hidden units followed by ReLU activations. Since the actions are to be scaled by the proportional gain constants in the PD controller, we apply Tanh activation to the controller network's output layer to ensure the values are normalized between -1 and +1. The value function outputs a scalar value

approximating the expected return which is not normalized. The input to the networks are as described earlier (§3.5).

## 4.3 Learning Full Model Policies Using Simplified Model Imitation

The full model policy is trained using imitation learning by tracking the reference trajectory generated by the simplified model. The imitation learning task can be considered as an inverse dynamics problem where the full model policy is required to produce the joint torques at each time step given the current character state and the future reference trajectory. We can consider the reward function to have a task reward component, an auxiliary reward component, and a style reward component.

As in the simplified environment, the task reward, $r_{\text{task}}$, is the sparse reward for successfully grabbing the next target handhold. The auxiliary reward ($r_{\text{aux}}$) consists of reward objectives that facilitate learning, including a tracking term and a reaching term. The tracking term rewards the character to closely follow the reference trajectory and a reaching term encourage the grabbing hand to be close to the next target handhold. The style reward ($r_{\text{style}}$) terms are added to incentivize more natural motions, including keeping the torso upright, slower arm rotation, reducing lower body movements, and minimizing energy. Finally, the overall reward can be computed as:

$$r_{\text{full}} = \exp(r_{\text{aux}} + r_{\text{style}}) + r_{\text{task}} \tag{2}$$

$$r_{\text{aux}} = w_t r_{\text{tracking}} + w_r r_{\text{reaching}} \tag{3}$$

$$r_{\text{style}} = w_u r_{\text{upright}} + w_a r_{\text{arm}} + w_l r_{\text{legs}} + w_e r_{\text{energy}} \tag{4}$$

Note that the style reward does not affect the learning outcome; the character can learn brachiation motions without the style reward. The reward terms and their weighting coefficients are fully described in Appendix A.

*Full Policy Networks.* The controller and the value function of the full model consist of two five-layer neural networks, with a layer width of 256. The first three hidden layers of the policy use the softsign activiation, while the final two hidden layers use ReLU activation. A Tanh is applied to the final output to normalize the actions. The critic uses ReLU for all hidden layers. Prior work is contradictory in nature with respect to the impact of network size, finding larger network sizes to perform both better [Wang et al. 2020] and worse [Bergamin et al. 2019] for full-body motion imitation tasks. In our experiments, we find the smaller networks used for the simplified model to be incapable of learning the brachiation task in the full environment.

## 4.4 Initial State and Early Termination

For both environments, we initialize the character to be hanging from the first handhold at a 90-degree angle, where zero degree corresponds to the pendulum resting position. The initial velocity of the character is set to zero as the base state provides the necessary forward momentum to reach the next handhold.

The environment resets when the termination condition is satisfied, which we define as entering an unrecoverable state. An unrecoverable state is reached when the character is not grabbing
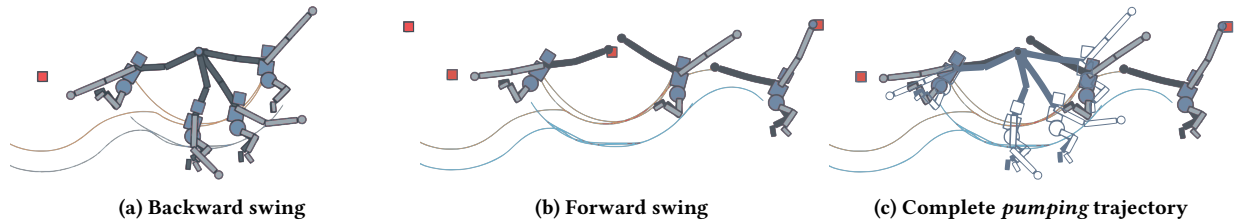
(a) Backward swing        (b) Forward swing        (c) Complete *pumping* trajectory

**Figure 3: Brachiation with emergent *pumping behavior*. The blue line gives the trajectory of the simplified model, while the brown line gives the trajectory of the root of the full character. The disparity between the two is caused by the fact that the simplified model has a longer arm length with respect to the full model, to better represent its center of mass. For further visual trajectories please refer to Appendix D and to the supplementary video.**

onto any handholds and falls below the graspable range of the next handhold with a downward velocity.

In addition to the unrecoverable criteria, both environments implement a minimum and maximum grab duration. The minimum and maximum duration are set to 0.25 and 4 seconds respectively. The idea of the minimum grab duration is to simulate the reaction time and time required to physically form a fist during a grasp. The maximum grab duration effectively serves as another form of early termination. It forces the character to enter an unrecoverable state unless the policy has already learned to generate forward momentum. We check the importance of early termination in Section 5.

## 5 RESULTS

All experiments are performed on a single 12-core machine with one NVIDIA RTX 2070 GPU. For the simplified model, training takes approximately five minutes as the simulation is parallelized on the GPU. Training the full model policy takes just under one hour. Each experiment is limited to collect a total of 25 million samples and the learning rate is decreased to $3 \times 10^{-5}$ over time using exponential annealing. We use the Adam [Kingma and Ba 2014] optimizer with a mini-batch size of 2000 and a learning rate of $3 \times 10^{-4}$ to update policy parameters in all experiments. Other implementation and hyperparameter details are summarized in Appendix C.

### 5.1 Simplified Model Results

The goal of the simplified model is to facilitate more efficient training of a full model policy by generating physically-feasible reference trajectories. A baseline solution for the simplified model can be successfully trained to find reasonable solutions for traversing difficult terrains, i.e. with handholds sampled from the full distribution (§3.3), using only the task reward described in Section 4.2. The learned controllers have some speed variations; we use a low speed controller to train the full model policy as the generated trajectories are closer to the motions demonstrated by real gibbons.

The learned control policies can traverse challenging sequences of handholds which require significant flight phases and exhibit emergent *pumping* behavior where the characters would perform extra back-and-forth swings as needed. We observe three scenarios where pumping behavior can be observed: waiting for the minimum swing time to elapse, adjusting for a better release angle, and for gaining momentum. A visualization of the pumping behavior is

shown in Figure 3. Please refer to the supplementary video for visual demonstrations of generated trajectories.

*Number of Look-ahead Handholds.* In this experiment, we empirically verify the choice for the number of look-ahead handholds by comparing the task performance on $N = \{1, 2, 3, 5, 10\}$. Results are summarized in Table 2. The best performance is achieved when $N = 1$ as measured by both the number of completed handholds and the average episode reward. This is surprising since previous work in human locomotion has shown that a two-step anticipation horizon achieves the lowest stepping error [Coros et al. 2008]. The result shows that the policy performance generally decreases with increasing number of look-ahead handholds. We hypothesize that, due to the relative small network size used in the simplified policy, increasing the number of look-ahead handholds creates more distraction for the policy. We leave validation of this hypothesis for future work.

*Importance of Early Termination.* In this ablation study, we verify the importance of applying early termination in the simplified environment. There are two forms of early termination: the unrecoverable criteria and the maximum grab duration. The complete simplified model implements all early termination strategies. Table 3 shows the average number of handholds reached aggregated across five different runs. The results show that enforcing strict early termination can facilitate faster learning. When maximum grab duration is disabled, only one run out of ten total successfully learned the brachiation task. In addition, we find that enforcing the minimum grab duration of 0.25 seconds helps the character to learn to swing, which in turn helps to discover the future handholds.

### 5.2 Full Model Results

The full model policy is learned through simplified model imitation, as described earlier (§4.3). In addition to using the tracking reward, which uses the simplified model only during training, we experiment with three different methods of further leveraging the simplified model reference trajectories at inference time.

**A. Tracking reward only**. The tracking reward is combined with the task reward to train the full model policy. The reference trajectory is only used to compute the tracking reward.

**B. Tracking + Other Rewards**. Reward includes all reward terms: tracking, task, auxiliary, and style rewards.

**Table 2: Experiment results on the effect of the number of look-ahead handholds for the simplified model. We report the average and the standard deviation over five independent runs.**

| | | Number of Look-ahead Handholds | | |
|---|---|---|---|---|
| 1 | 2 | 3 | 5 | 10 |
| *Handholds Completed* | | | | |
| **17.8 ± 2.1** | 14.6 ± 5.0 | 4.9 ± 5.7 | 7.6 ± 5.8 | 0.3 ± 0.6 |
| *Episode Reward (×10²)* | | | | |
| **7.3 ± 0.8** | 5.9 ± 1.8 | 2.0 ± 2.4 | 3.1 ± 2.4 | 0.1 ± 0.3 |

**Table 3: Ablation results on early termination strategies for the simplified model. Each column names the termination strategy that is removed. Each entry represents the average number of completed handholds and the standard deviation over five runs.**

| | Early Termination Ablations | | | |
|---|---|---|---|---|
| None | Unrecoverable | Min Duration | Max Duration | Min & Max |
| **14.6 ± 5.1** | 10.2 ± 7.3 | 7.4 ± 9.3 | 2.9 ± 2.0 | 3.1 ± 4.1 |

**C. Rewards + Release Timing**. Reward includes everything described in Experiment B. In addition, the release timing from the reference trajectory is used to control the grab action of the full model policy.

**D. Rewards + States**. Reward includes everything described in Experiment B. In addition, a slice of the future reference trajectory is included as a part of the full model policy input.

**E. Rewards + States + Grab Info**. Reward includes everything described in Experiment D. In addition, instead of conditioning the full model state with just the future reference trajectory, we also give access to the grab flag for those future points.

Figure 4 shows the learning curves for different full model configurations. We include a baseline model which is trained using all reward terms except for the tracking reward. For the baseline model, we experimented both with and without a learning curriculum on the terrain difficulty, e.g., similar to [Xie et al. 2020], but neither version was able to advance to the second handhold. The best policies are obtained when the simplified model is also used at evaluation time. This includes adding information from the reference trajectory (Experiment D) and also further adding the grab information (Experiment E). However, the grab information provides no additional benefit to the policy when reference trajectory information is present.

## 5.3　Full Model with Planning

At inference time, we are able to extend the capability of the full model policy by expanding the anticipation horizon using the simplified model. The full model policy's value function cannot anticipate more than the number of look-ahead handholds that was used for training. However, the simplified model is fast enough to replan
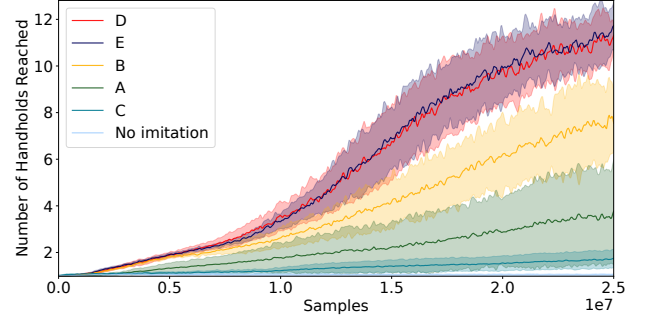


**Figure 4: Learning Curves for the full model, for various configurations. No imitation reward results in failure. Using reference trajectory information only at training time (A and B) improves learning. Using the release timing from the simplified model degrades performance (C). The best policies are obtained when simplified and full models are used in tandem at evaluation time (D and E). Please refer to Section 5.2 for a detailed description.**

**Table 4: Quantitative evaluation of the different planning methods. In each generated terrain, three unreachable gaps (e.g., Figure 5) are placed at random locations. Each planning method is evaluated on 40 different terrains. The combined approach is able to pass all gaps most consistently.**
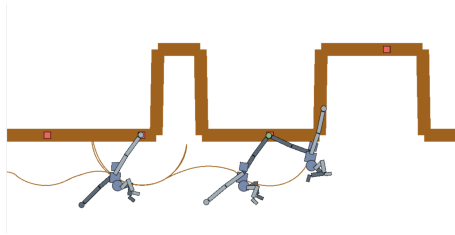
| | Gaps Passed | | | |
|---|---|---|---|---|
| Planning Method | 0 | 1 | 1+2 | 1+2+3 |
| Only Full Value Function | 40 | 22 | 9 | 7 |
| Only Simplified Rewards | 40 | 36 | 26 | 15 |
| Reward & Value Function | 40 | 35 | 26 | **21** |

on a handhold-by-handhold basis, in the style of model predictive control (MPC). It can be used to simulated thousands of trajectories in parallel at run-time. In this experiment, we set the planner to explore 10k trajectories for 4 seconds into the future; this can be completed in under one second of machine time. There is still room for improvement for real time animation purposes.
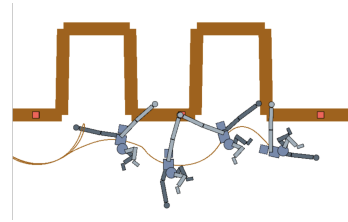
The planner uses three fully-trained components: the full model policy $\pi_{\text{full}}$, the full model value function $V_{\text{full}}$ and the simplified model policy $\pi_{\text{simple}}$. As input, the method takes a terrain which is a contiguous ceiling where handholds can be placed. Given the current state of the full model character $s_t$, we randomly sample $K$ handhold plans of length $H$ handholds on the terrain, $\{\mathcal{P}_k\}$. Concretely, each handhold in a plan $\mathcal{P}_k$ is sampled at a horizontal distance $\delta x \sim U(1.1, 1.8)$ beyond the previous handhold, and takes its corresponding height from the terrain. We then retain the best plan, i.e., $k = \arg \max J(\mathcal{P}_k)$, according to the objective defined by

$$J(\mathcal{P}_k) = V_{\text{full}}(s_t, k[0:N]) + \sum_j^H R_j,$$

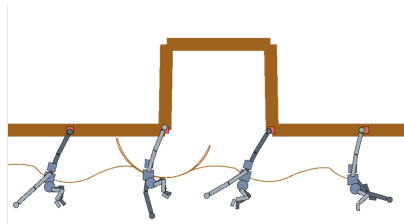where $R_j$ is the reward for the simplified model obtained by $\pi_{\text{simple}}$, and $N$ is the policy look-ahead value, with $H > N$. The selected plan is then used until the next handhold is reached, at which point

(a) Planning using only the full model value function. The last planned target grasp location (red dot) is unreachable. This approach has a limited planning horizon, so unrealistic plans are occasionally generated.



(b) Planning using only the simplified model. Here, the last handhold target is unreachable because this planning approach fails to consider the capability of the full model policy.



(c) Planning with combined full model value function and simplified model. The target handhold locations are planned very close to the discontinuities, which effectively reduces the risk of producing unreachable targets.

**Figure 5: Comparison of different planning approaches.**

a full replanning takes place. This MPC-like process repeats *ad infinitum*. We find empirically that the above objective generates better planning trajectories as compared to only using the reward of $\pi_{\text{simple}}$ or only $V_{\text{full}}$. This is because the reward of the simplified model is not enough to capture the capabilities of the full character, while the value function has a limited planning horizon. Quantitative results are presented in Table 4. Examples of terrains traversed by the policy with the planners are presented in Figure 5 and in the supplementary video.

## ACKNOWLEDGMENTS

## 6  CONCLUSIONS

Brachiation is challenging to learn, as it requires the careful management of momentum, as well as precision in grasping. We demonstrate a two-level reinforcement learning method for learning highly-capable and dynamic brachiation motions. The motions generated by the policy learned for the simplified model play a critical role in the efficient learning of the full policy. Our work still has many limitations. There are likely to be other combinations of rewards, reward curricula, and handhold curricula, that, when taken together, may produce equivalent capabilities. We have not yet fully characterized the limits of the model's capabilities, i.e., the simulated gibbon's ability to climb, leap, descend, and more. We wish to better understand, in a quantitative fashion, how the resulting motions compare to those of real gibbons. Lastly, to emulate the true capabilities of gibbons, we need methods that can plan in complex 3D worlds, perform 3D brachiation, make efficient use of the legs when necessary, and more. We believe the simplicity of our approach provides a solid foundation for future research in these directions.

# REFERENCES

Kevin Bergamin, Simon Clavet, Daniel Holden, and James Richard Forbes. 2019. DReCon: Data-Driven Responsive Control of Physics-Based Characters. *ACM Trans. Graph.* 38, 6, Article 206 (nov 2019), 11 pages. https://doi.org/10.1145/3355089.3356536

Glen Berseth. 2019. *Scalable deep reinforcement learning for physics-based motion control.* Ph. D. Dissertation. University of British Columbia. https://doi.org/10.14288/1.0378079

J. E. Bertram, A. Ruina, C. E. Cannon, Y. H. Chang, and M. J. Coleman. 1999. A point-mass model of gibbon locomotion. *The Journal of Experimental Biology* 202, Pt 19 (1999), 2609–2617.

Nuttapong Chentanez, Matthias Müller, Miles Macklin, Viktor Makoviychuk, and Stefan Jeschke. 2018. Physics-based motion capture imitation with deep reinforcement learning. In *Proceedings of the 11th Annual International Conference on Motion, Interaction, and Games.* 1–10.

Stelian Coros, Philippe Beaudoin, Kang Kang Yin, and Michiel van de Panne. 2008. Synthesis of Constrained Walking Skills. In *ACM SIGGRAPH Asia 2008 Papers* (Singapore) (SIGGRAPH Asia '08). Association for Computing Machinery, New York, NY, USA, Article 113, 9 pages. https://doi.org/10.1145/1457515.1409066

Erwin Coumans and Yunfei Bai. 2016–2022. PyBullet, a Python module for physics simulation for games, robotics and machine learning. http://pybullet.org.

Siavash Farzan, Ai-Ping Hu, Evan Davies, and Jonathan Rogers. 2018. Modeling and Control of Brachiating Robots Traversing Flexible Cables. In *2018 IEEE International Conference on Robotics and Automation (ICRA).* 1645–1652. https://doi.org/10.1109/ICRA.2018.8461036 ISSN: 2577-087X.

John Fleagle. 1974. Dynamics of a brachiating siamang [Hylobates (Symphalangus) syndactylus]. *Nature* 248, 5445 (1974), 259–260. https://doi.org/10.1038/248259a0 Bandiera_abtest: a Cg_type: Nature Research Journals Number: 5445 Primary_atype: Research Publisher: Nature Publishing Group.

Siddhant Gangapurwala, Mathieu Geisert, Romeo Orsolino, Maurice Fallon, and Ioannis Havoutis. 2020. Rloc: Terrain-aware legged locomotion using reinforcement learning and optimal control. *arXiv preprint arXiv:2012.03094* (2020).

Siddhant Gangapurwala, Mathieu Geisert, Romeo Orsolino, Maurice Fallon, and Ioannis Havoutis. 2021. Real-Time Trajectory Adaptation for Quadrupedal Locomotion using Deep Reinforcement Learning. In *Proceedings of the International Conference on Robotics and Automation (ICRA).* Institute of Electrical and Electronics Engineers.

Mario W. Gomes and Andy L. Ruina. 2005. A five-link 2D brachiating ape model with life-like zero-energy-cost motions. *Journal of Theoretical Biology* 237, 3 (2005), 265–278. https://doi.org/10.1016/j.jtbi.2005.04.014

Kevin Green, Yesh Godse, Jeremy Dao, Ross L. Hatton, Alan Fern, and Jonathan Hurst. 2021. Learning Spring Mass Locomotion: Guiding Policies with a Reduced-Order Model. *arXiv:2010.11234 [cs]* (March 2021). http://arxiv.org/abs/2010.11234 arXiv: 2010.11234.

Hideki Kajima, Masahiro Doi, Yasuhisa Hasegawa, and Toshio Fukuda. 2004. A study on a brachiation controller for a multi-locomotion robot — realization of smooth, continuous brachiation. *Advanced Robotics* 18, 10 (2004), 1025–1038. https://doi.org/10.1163/1568553042674671 Publisher: Taylor & Francis _eprint: https://doi.org/10.1163/1568553042674671.

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

Ilya Kostrikov. 2018. PyTorch Implementations of Reinforcement Learning Algorithms. https://github.com/ikostrikov/pytorch-a2c-ppo-acktr-gail.

Scott Kuindersma. 2020. Recent Progress on Atlas, the World's Most Dynamic Humanoid Robot. https://www.youtube.com/watch?v=EGABAx52GKI

Taesoo Kwon, Yoonsang Lee, and Michiel van de Panne. 2020. Fast and Flexible Multilegged Locomotion Using Learned Centroidal Dynamics. *ACM Trans. Graph.* 39, 4, Article 46 (jul 2020), 24 pages. https://doi.org/10.1145/3386569.3392432

Chi-Ying Lin and Zong-Han Yang. 2020. TRBR: Flight body posture compensation for transverse ricochetal brachiation robot. *Mechatronics* 65 (2020), 102307. https://doi.org/10.1016/j.mechatronics.2019.102307

Ying-Sheng Luo, Jonathan Hans Soeseno, Trista Pei-Chun Chen, and Wei-Chao Chen. 2020. CARL: Controllable Agent with Reinforcement Learning for Quadruped Locomotion. *ACM Trans. Graph.* 39, 4, Article 38 (jul 2020), 10 pages. https://doi.org/10.1145/3386569.3392433

Josh Merel, Leonard Hasenclever, Alexandre Galashov, Arun Ahuja, Vu Pham, Greg Wayne, Yee Whye Teh, and Nicolas Heess. 2018. Neural probabilistic motor primitives for humanoid control. *arXiv preprint arXiv:1811.11711* (2018).

Fana Michilsens, Evie E Vereecke, Kristiaan D'Août, and Peter Aerts. 2009. Functional anatomy of the gibbon forelimb: adaptations to a brachiating lifestyle. *Journal of Anatomy* 215, 3 (2009), 335–354. https://doi.org/10.1111/j.1469-7580.2009.01109.x

Sehee Min, Jungdam Won, Seunghwan Lee, Jungnam Park, and Jehee Lee. 2019. Soft-Con: Simulation and Control of Soft-Bodied Animals with Biomimetic Actuators. *ACM Trans. Graph.* 38, 6, Article 208 (2019).

J. Nakanishi, T. Fukuda, and D.E. Koditschek. 2000. A brachiating robot controller. *IEEE Transactions on Robotics and Automation* 16, 2 (2000), 109–123. https://doi.org/10.1109/70.843166 Conference Name: IEEE Transactions on Robotics and Automation.

Soohwan Park, Hoseok Ryu, Seyoung Lee, Sunmin Lee, and Jehee Lee. 2019. Learning predict-and-simulate policies from unorganized human motion data. *ACM Transactions on Graphics (TOG)* 38, 6 (2019), 1–11.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. PyTorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems.* 8024–8035.

Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. 2018. Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Transactions on Graphics (TOG)* 37, 4 (2018), 1–14.

Ioannis Poulakakis and Jessy W Grizzle. 2009. The spring loaded inverted pendulum as the hybrid zero dynamics of an asymmetric hopper. *IEEE Trans. Automat. Control* 54, 8 (2009), 1779–1793.

Holger Preuschoft and Brigitte Demes. 1985. Influence of Size and Proportions on the Biomechanics of Brachiation. In *Size and Scaling in Primate Biology*, William L. Jungers (Ed.). Springer US, Boston, MA, 383–399. https://doi.org/10.1007/978-1-4899-3647-9_17

Fuminori Saito. 1992. Movement control of brachiation robot using CMAC between different distance and height. In *Proc. IMACS/SICE Int. Symp. on Robotics, Mechatronics and Manufacturing Systems.* 35–40.

F. Saito, T. Fukuda, and F. Arai. 1994. Swing and locomotion control for a two-link brachiation robot. *IEEE Control Systems Magazine* 14, 1 (1994), 5–12. https://doi.org/10.1109/37.257888 Conference Name: IEEE Control Systems Magazine.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* (2017).

Bruno Siciliano, Oussama Khatib, and Torsten Kröger. 2008. *Springer handbook of robotics.* Vol. 200. Springer.

Nitish Sontakke and Sehoon Ha. 2021. Solving Challenging Control Problems Using Two-Staged Deep Reinforcement Learning. *arXiv preprint arXiv:2109.13338* (2021).

Vassilios Tsounis, Mitja Alge, Joonho Lee, Farbod Farshidian, and Marco Hutter. 2020. DeepGait: Planning and Control of Quadrupedal Gaits using Deep Reinforcement Learning. *arXiv:1909.08399 [cs]* (Jan. 2020). http://arxiv.org/abs/1909.08399 arXiv: 1909.08399.

Dengke Wan, Hongtai Cheng, Guangfei Ji, and Shuai Wang. 2015. Non-horizontal ricochetal brachiation motion planning and control for two-link Bio-primate robot. In *2015 IEEE International Conference on Robotics and Biomimetics (ROBIO).* 19–24. https://doi.org/10.1109/ROBIO.2015.7407033

Tingwu Wang, Yunrong Guo, Maria Shugrina, and Sanja Fidler. 2020. UniCon: Universal Neural Controller For Physics-based Character Motion. arXiv:2011.15119 [cs.GR]

Jungdam Won, Deepak Gopinath, and Jessica Hodgins. 2020. A scalable approach to control diverse behaviors for physically simulated characters. *ACM Transactions on Graphics (TOG)* 39, 4 (2020), 33–1.

Jungdam Won, Jungnam Park, and Jehee Lee. 2018. Aerobatics Control of Flying Creatures via Self-Regulated Learning. *ACM Trans. Graph.* 37, 6, Article 181 (dec 2018), 10 pages. https://doi.org/10.1145/3272127.3275023

Zhaoming Xie, Xingye Da, Buck Babich, Animesh Garg, and Michiel van de Panne. 2021. GLiDE: Generalizable Quadrupedal Locomotion in Diverse Environments with a Centroidal Model. *arXiv:2104.09771 [cs]* (April 2021). http://arxiv.org/abs/2104.09771 arXiv: 2104.09771.

Zhaoming Xie, Hung Yu Ling, Nam Hee Kim, and Michiel van de Panne. 2020. ALLSTEPS: Curriculum-driven Learning of Stepping Stone Skills. *Computer Graphics Forum (Proc. ACM SIGGRAPH/EG Symposium on Computer Animation)* 39, 8 (2020), 213–224.

Shuo Yang, Zhaoyuan Gu, Ruohai Ge, Aaron M. Johnson, Matthew Travers, and Howie Choset. 2019. Design and Implementation of a Three-Link Brachiation Robot with Optimal Control Based Trajectory Tracking Controller. *arXiv:1911.05168 [cs, eess]* (2019). http://arxiv.org/abs/1911.05168 arXiv: 1911.05168.

Zheng Zhang and Kok Cheong Wong. 1999. Details and Implementation Issues of Animating Brachiation. In *Computer Animation and Simulation '99 (Eurographics)*, Nadia Magnenat-Thalmann and Daniel Thalmann (Eds.). Springer, Vienna, 123–132. https://doi.org/10.1007/978-3-7091-6423-5_12

## A  FULL MODEL REWARDS

The tracking reward $r_{\text{tracking}}$ is formulated as a penalty and computed as the squared distance between the character's root and the position of the reference trajectory.

$$r_{\text{tracking}} = \|p_{\text{body}} - p_{\text{reference}}\|_2^2$$
$$w_t = -4$$

The reaching reward $r_{\text{reaching}}$ is computed similarly as the tracking reward, except the distance is between the next grabbing hand and the target handhold. This reward is only applied during the flight phase. During swing phase, the character can still use the free arm to generate momentum without penalty.

$$r_{\text{reaching}} = \|p_{\text{hand}} - p_{\text{target}}\|_2^2, \text{ if in flight phase}$$
$$w_r = -0.1$$

The upright posture term $r_{\text{upright}}$ penalizes the character if the root pitch is not within 40 degrees of the vertical axis.

$$r_{\text{upright}} = |\text{pitch}| - 40°, \text{ if } |\text{pitch}| > 40°$$
$$w_u = -1$$

The arm rotation reward $r_{\text{arm}}$ penalizes the character for excessive angular velocity in the next grabbing arm.

$$r_{\text{arm}} = |\omega_{\text{arm}}| \text{ and } w_a = -0.1$$

The lower body reward $r_{\text{legs}}$ encourages the knees to be close to 110 degrees from base position of straight legs. This reward is applied to improve the visibility of the legs during brachiation.

$$r_{\text{legs}} = \|\theta_{\text{knees}} - 110°\|_1$$
$$w_l = -0.1$$

The energy reward $r_{\text{energy}}$ penalizes the policy for using excessive amounts of energy. The overall energy is approximated based on the applied torques and the combined angular velocity of all joints, computed as:

$$r_{\text{energy}} = \|\tau_{\text{joints}}\|_2^2 + \|\omega_{\text{joints}}\|_1$$
$$w_l = -0.01$$

## B  PROXIMAL POLICY OPTIMIZATION

Let an experience tuple be $e_t = (o_t, a_t, o_{t+1}, r_t)$ and a trajectory be $\tau = \{e_0, e_1, \ldots, e_T\}$. We episodically collect trajectories for a fixed number of environment transitions and we use this data to train the controller and the value function networks. The value function network approximates the expected future returns of each state, and is defined for a policy $\pi$ as

$$V^\pi(o) = E_{o_0=o, a_t \sim \pi(\cdot|o_t)} \left[ \sum_{t=0}^{\infty} \gamma^t r(o_t, a_t) \right].$$

This function can be optimized using supervised learning due to its recursive nature:

$$V^{\pi_\theta}(o_t) = \gamma V^{\pi_\theta}(o_{t+1}) + r_t,$$

where

$$V^{\pi_\theta}(o_T) = r_T + \gamma V^{\pi_{\theta_{old}}}(o_{T+1}).$$

In PPO, the value function is used for computing the advantage

$$A_t = V^{\pi_\theta} - V^{\pi_{\theta_{old}}}$$

which is then used for training the policy by maximizing:

$$L_\pi(\theta) = \frac{1}{T} \sum_{t=1}^{T} \min(\rho_t \hat{A}_t, \text{ clip}(\rho_t, 1-\epsilon, 1+\epsilon)\hat{A}_t),$$

where $\rho_t = \pi_\theta(a_t|o_t) / \pi_{\theta_{old}}(a_t|o_t)$ is an importance sampling term used for calculating the expectation under the old policy $\pi_{\theta_{old}}$.

## C  PPO HYPERPARAMETERS

**Table 5: Hyperparameters used for training PPO.**

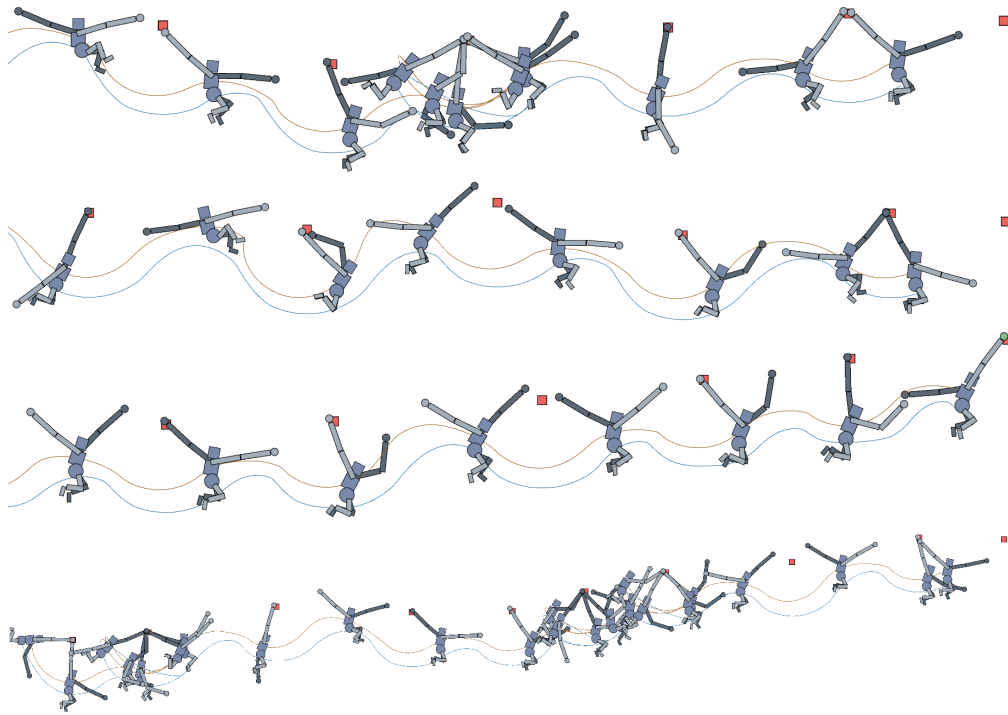| Hyperparameter | Simplified/Full Model |
|---|---|
| Learning Rate | $3 \times 10^{-4}$ |
| Final Learning Rate | $3 \times 10^{-5}$ |
| Optimizer | Adam |
| Batch Size | 2000 |
| Training Steps | $2.5 \times 10^7$ |
| Num processes | 10000/125 |
| Episode Steps | 80000/40000 |
| Num PPO Epochs | 10 |
| Discount Factor $\gamma$ | 0.99 |
| Gradient Clipping | False |
| Entropy Coefficient | 0 |
| Value Loss Coefficient | 0 |
| Clip parameter | 0.2 |
| Max Grad Norm | 2.0 |

# D SUPPLEMENTARY RESULTS



**Figure 6: Additional trajectories showing a variety of motions.**