
Improving the Representation Learning Ability of PPO Algorithm by Interleaving States and Actions

Tianze Luo
IIIS, Tsinghua University
2022012084

Xingchen Miao
IIIS, Tsinghua University
2022012203

Abstract

In this work, we significantly improve the data efficiency and performance of the Proximal Policy Optimization (PPO) algorithm on HumanoidBench, a benchmark for robot learning in high-dimensional simulated environments. HumanoidBench features a humanoid robot with dexterous hands and presents a variety of challenging tasks involving both whole-body manipulation and locomotion. According to the original HumanoidBench paper, the PPO algorithm fails to learn any practical policy within 10 million sample steps, even on relatively simple tasks like walking. To enhance the sample efficiency and performance of PPO, we improve its representation learning capability by sequentially modeling the policy-value network with interleaved states and actions as inputs. We also implement several beneficial techniques such as state normalization, reward scaling, Generalized Advantage Estimation (GAE), and carefully tune the hyperparameters, including loss function scales and the training schedule. Our improved PPO algorithm successfully learns effective policies and significantly outperforms the model-free methods reported in the original HumanoidBench paper, including the baseline PPO and Soft Actor-Critic (SAC) algorithms. Moreover, our method achieves performance similar to model-based approaches like DreamerV3 and TD-MPC2, demonstrating its potential for solving complex tasks in humanoid robot learning.

1 Introduction

Humanoid robots, with their human-like morphology, hold great potential in a wide range of applications, from assisting with daily human tasks to functioning autonomously in dynamic environments. However, the development of effective controllers for humanoid robots has been hampered by the high costs and fragility of real-world hardware setups. To mitigate these challenges, HumanoidBench was introduced as a high-dimensional simulated benchmark for robot learning, designed to test algorithms in environments that simulate complex humanoid robot tasks. Built on the MuJoCo physics engine, HumanoidBench features a humanoid robot (Unitree H1) equipped with dexterous hands, and supports a variety of locomotion and manipulation tasks. These tasks span from basic walking and reaching to more intricate whole-body manipulations like unloading packages, cleaning windows, and playing basketball. With these features, HumanoidBench serves as both a testbed for evaluating robot learning algorithms and a platform for advancing research in humanoid robotics, offering the research community a controlled, safe, and cost-effective means to benchmark algorithms in complex, long-horizon tasks. The challenges faced by algorithms on this platform include not only the complexity of robot dynamics but also the need for coordination across multiple body parts and the intricacies of high-dimensional action spaces.

Previous research on robotic simulation has predominantly concentrated on tasks characterized by simpler dynamics and shorter time horizons. These tasks are often limited to single-arm manipulation or relatively straightforward locomotion problems. While benchmarks such as Robosuite and other

bimanual manipulation environments have made significant progress in addressing coordination challenges, they still fall short of capturing the intricate full-body coordination required by humanoid robots. Some benchmarks have incorporated long-horizon tasks, like block stacking or furniture assembly, but these remain constrained in their ability to allow robots to engage with their environment as holistic, fully integrated systems. To address these complexities, hierarchical reinforcement learning (HRL) has been explored as a strategy to separate low-level control from high-level decision-making, which can aid in managing long-horizon tasks. However, scaling HRL to humanoid robots presents additional challenges, as low-level RL typically requires extensive parallel pretraining, while high-level RL demands intricate coordination across a multitude of degrees of freedom.

One of the core challenges in humanoid robot learning lies in the efficient utilization of data. The high-dimensional nature of humanoid systems, coupled with the complexity of their environments, often results in learning processes that are both slow and data-intensive. Model-based approaches, such as TD-MPC2 and DreamerV3, have demonstrated progress in addressing these issues by learning and leveraging world models. In contrast, model-free methods, such as Soft Actor-Critic (SAC) and Proximal Policy Optimization (PPO), are constrained by their inability to efficiently reuse past data, as they lack a learned world model and rely solely on real samples collected during training.

The limitations of PPO are particularly pronounced due to its on-policy nature, which prevents the algorithm from repeatedly reusing data stored in a replay buffer—a capability that off-policy methods like SAC employ effectively. This inability to make multiple passes over the same data exacerbates the challenges of scaling PPO to complex tasks. On benchmarks such as HumanoidBench, which demand both precise high-dimensional control and robust long-horizon decision-making, these deficiencies render PPO insufficient for achieving practical performance. Indeed, PPO often fails even on relatively simple tasks, such as basic walking, further highlighting its limitations in addressing the unique challenges posed by humanoid robots.

In this work, we address the shortcomings of PPO on HumanoidBench by improving its data efficiency and performance. We propose a novel approach to PPO by incorporating a sequential policy-value network, where states and actions are interleaved as inputs, thereby enhancing the network’s representation learning capability. This approach allows the model to capture the complex dependencies between states and actions in long-horizon tasks, improving both stability and efficiency in learning. By tuning key hyperparameters and integrating techniques such as state normalization, reward scaling, and Generalized Advantage Estimation (GAE), our method significantly outperforms existing model-free methods including SAC and origin PPO, which achieves similar performances to state-of-the-art model-based algorithms like DreamerV3 and TD-MPC2.

2 Related Work

Proximal Policy Optimization (PPO) is a widely adopted reinforcement learning (RL) algorithm, renowned for its simplicity and effectiveness across a wide range of tasks. PPO improves upon traditional policy gradient methods by introducing a surrogate objective function that constrains policy updates to remain within a stable range. The key innovation of PPO lies in its "clipped" objective, which prevents the new policy from deviating excessively from the old one during updates. By bounding the ratio of new to old action probabilities, PPO mitigates instability often associated with large policy shifts. Despite its widespread success, PPO’s on-policy, model-free nature limits its ability to reuse past experiences, resulting in inefficiencies when tackling complex, high-dimensional tasks, such as those found in the HumanoidBench benchmark. These limitations become particularly problematic in scenarios where sample efficiency is crucial. The pseudocode for the PPO algorithm is described in Algorithm 1.

PPO has several strengths that have contributed to its popularity. Its clipped objective ensures stable training by constraining updates, making it particularly effective in environments with high reward variance. Moreover, PPO is relatively easy to implement compared to alternatives like Trust Region Policy Optimization (TRPO), which requires a more complex constraint handling mechanism. PPO’s versatility has also been demonstrated across a wide variety of tasks, ranging from robotic control to video games. However, PPO is not without its limitations. Its on-policy nature requires fresh data at every iteration, which limits the reuse of past experiences and leads to high sample complexity. Additionally, PPO is highly sensitive to hyperparameter choices, such as clipping thresholds, learning rates, and parameters for advantage estimation, which can significantly affect performance if not

Algorithm 1 Proximal Policy Optimization (PPO)

```
1: Initialize: Policy  $\pi$ , value function  $V$ , and hyperparameters such as clipping threshold  $\epsilon$ .
2: for each iteration do
3:   Collect Rollouts: Interact with the environment using policy  $\pi$ , collecting  $N$  samples of
       $(s_t, a_t, r_t)$ .
4:   Compute Returns and Advantages:
5:     for each sample  $t$  in the rollout do
6:       Compute reward-to-go:  $R_t = \sum_{t'=t}^T \gamma^{t'-t} r_{t'}$ .
7:       Estimate advantage  $A_t$  using Generalized Advantage Estimation (GAE):
          
$$A_t = \delta_t + (\gamma\lambda)\delta_{t+1} + \dots + (\gamma\lambda)^{T-t+1}\delta_{T-1},$$

      where  $\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$ .
8:     end for
9:     for each minibatch in the sampled data do
10:      Compute new action log probabilities:  $\log \pi(a_t|s_t)$ .
11:      Compute probability ratio:
          
$$\text{ratio} = \frac{\pi_{\text{new}}(a_t|s_t)}{\pi_{\text{old}}(a_t|s_t)} = \exp(\log \pi_{\text{new}} - \log \pi_{\text{old}}).$$

12:      Clip the probability ratio:
          
$$\text{clipped\_ratio} = \min(\max(\text{ratio}, 1 - \epsilon), 1 + \epsilon).$$

13:      Calculate policy loss:
          
$$\text{policy\_loss} = -\mathbb{E}_t [\min(\text{ratio} \cdot A_t, \text{clipped\_ratio} \cdot A_t)].$$

14:      Calculate value loss:
          
$$\text{value\_loss} = \mathbb{E}_t [(R_t - V(s_t))^2].$$

15:      Optionally, calculate entropy loss:
          
$$\text{entropy\_loss} = -\mathbb{E}_t [\pi(a_t|s_t) \cdot \log \pi(a_t|s_t)].$$

16:      Combine losses into a total objective:
          
$$\text{total\_loss} = \text{policy\_loss} + c_1 \cdot \text{value\_loss} - c_2 \cdot \text{entropy\_loss}.$$

17:      Update  $\pi$  and  $V$  via backpropagation using the total loss.
18:     end for
19: end for
```

carefully tuned. Furthermore, PPO struggles with tasks requiring long-term credit assignment due to its reliance on shorter-horizon advantage estimates.

3 Methods

3.1 Overview

PPO is a well-known reinforcement learning algorithm that has achieved significant success across a variety of tasks. However, on the HumanoidBench benchmark, it struggles to learn any practical policy. We identify two key factors contributing to this challenge. First, PPO is a model-free, on-policy algorithm, which limits its ability to leverage data efficiently. Unlike off-policy methods such as Soft Actor-Critic (SAC), PPO cannot take advantage of trajectory data from a world model, nor can it fully utilize previously sampled data, thereby hindering its performance. Second, the complexity of the HumanoidBench environment exacerbates this issue. The H1 robot in HumanoidBench has an action space with 61 dimensions, while the observations consist of hundreds of dimensions. Additionally, the action frequency is 50Hz, meaning that the relationships between states and actions are long-term, highly complex, and difficult to capture. This makes learning an effective policy with

PPO even more challenging, as it requires the model to handle and make decisions based on this intricate temporal and spatial information.

Based on these observations, we propose modeling the policy-value network in a sequential manner, where recent states and actions are interleaved as inputs. This design choice is inspired by the GLM-4-Voice, a recent large language model for speech-to-text tasks, which interleaves text and speech tokens during the pretraining phase to enhance the model’s multimodal understanding and generative capabilities. By borrowing this approach, we aim to capture the intricate dependencies between states and actions in a more effective and dynamic way. The structure of our policy-value network is shown below, and it stands in stark contrast to conventional architectures such as multi-layer perceptrons (MLPs) and recurrent neural networks (RNNs), offering a novel approach to sequential learning in reinforcement learning environments. To begin with, we utilize only the most recent 8

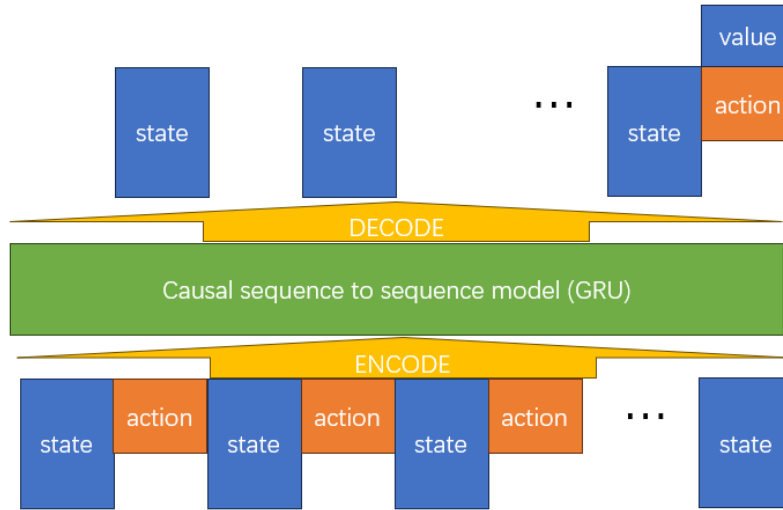


Figure 1: The sequential policy-value network

states and actions to predict the current action and value, striking a good balance between the full Markov prediction typically implemented by multi-layer perceptrons (MLPs) and the full sequential prediction used in recurrent neural networks (RNNs). In our model, both the encoder and decoder are implemented as MLPs, while the sequential modeling component is a two-layer Gated Recurrent Unit (GRU). The most recent state represents the current state from which we make predictions, and the actions are generated by an MLP utilizing the Beta distribution reparameterization trick. The reason we choose a linearly scaled Beta distribution over a clipped Gaussian distribution is that, in HumanoidBench, all actions are normalized to the range $[-1, 1]$. The Beta distribution naturally provides bounded outputs, making it a more suitable choice for this setting. The value function is predicted by a separate MLP, which takes the output from the final time step of the GRU as its input. This architecture efficiently combines the benefits of both Markovian and sequential dependencies for enhanced performance.

The core design principle of our approach is to enhance the representation learning capability of the network in order to improve both data efficiency and overall performance. However, we aim to avoid the world-model bias that can arise from inaccurate trajectory predictions, as seen in model-based methods. To mitigate this, we treat state prediction loss solely as an auxiliary loss, helping the model grasp the underlying physical laws of the environment, rather than using it to sample trajectories, which could introduce bias into the policy-value network. Additionally, it is worth noting that, unlike typical sequence modeling, our model does not predict actions. This is because, in reinforcement learning, actions are generated by the model itself, not from the data distribution, making action predictions an unreliable way to guide the network’s learning. We also experimented with using the average of previous action predictions to output the final state, but did not observe significant improvements compared to the approach we currently employ.

140 It is important to note that we combine the tasks of state, action, and value prediction into a single
141 network, which effectively enhances both data efficiency and the representation learning capability of
142 our model. However, this also introduces a challenge in balancing the associated loss functions. To
143 encourage exploration, we include an action entropy loss. Thus, the complete loss function for our
144 network consists of four components in total. After conducting several experiments, we empirically
145 set the coefficients for the state prediction loss, PPO-clip loss, value prediction loss, and action
146 entropy loss to 0.1, 1, 0.3, and 0.05, respectively, in order to stabilize the learning process.

147 In our view, our method offers three key advantages over the original PPO algorithm. First, our
148 model leverages the environment’s dynamics to assist in decision-making and value prediction, while
149 avoiding the world-model bias that can affect model-based approaches. Second, by using the 8
150 most recent states and actions to predict the current action, our model is better able to learn more
151 consistent actions over long-horizon tasks. Finally, our model achieves much higher data efficiency,
152 as all predictions are made by a unified network with stronger representation power (comparing to
153 MLPs). This significantly alleviates the challenges faced by the original PPO algorithm in terms of
154 data efficiency and learning stability.

155 Simply having a working algorithm is not sufficient; we also incorporate several auxiliary techniques
156 to stabilize the learning process. First, we use Generalized Advantage Estimation (GAE) with a
157 decay factor of 0.95 to estimate the advantages in the PPO-clip loss function, and the advantages
158 are normalized again in every batch. Additionally, the reward is scaled by a factor of 0.25, and the
159 discount factor is set to 0.99. Empirically, we found that the target for the value network should not
160 exceed 20, and the reward scaling coefficient of 0.25 is effective for tasks with rewards under 1000
161 over 1000 steps. For other cases, different scaling methods might be required. We also experimented
162 with adaptive reward normalization and scaling methods but observed worse performance compared
163 to simply applying a fixed scaling coefficient.

164 Additionally, due to the significant variation in environments across different tasks in HumanoidBench,
165 we employ state normalization to stabilize the input level of the network and reduce the impact of such
166 environmental differences. Specifically, we maintain a running estimate of the mean and standard
167 deviation of the states during the sampling process. These statistics are then used to normalize the
168 state data before feeding it into the network, ensuring that the input features remain within a consistent
169 range. To further improve training stability and convergence, we adopt orthogonal initialization for
170 the network weights, which helps maintain healthy gradients throughout training. Moreover, we
171 carefully tune the hyperparameters to stabilize the optimization process: the PPO-clip scale is set to
172 0.2, and the gradient clipping scale is set to 0.5.

173 All these choices mentioned above ensure that the learning process remains smooth and prevent
174 excessively large updates that could destabilize training.

175 Now, we will present some important specific implementation details of our algorithm to provide
176 a clearer understanding of its design and operation. The first implementation detail is the action
177 predictor. We predict the alpha and beta parameters using MLPs, with the Softplus function (+1.0)
178 applied to ensure that these parameters are strictly greater than one. Afterward, the Beta distribution
179 is scaled and shifted to the range $[-1, 1]$ to generate the actions. To prevent potential numerical
180 instability, we clamp the values at the boundaries.

181 Additionally, to accelerate the training process, we utilize 32 parallel processes for data sampling,
182 leveraging the SubprocVecEnv() for efficient parallelization. The total training consists of 10 million
183 sample steps, with each round involving the collection of 32 parallel environments, each sampling
184 1000 steps. This setup provides 32,000 samples per training period. During training, the batch size is
185 set to 500, and each round consists of 10 epochs, meaning the network is updated 640 times per round.
186 Given that our model requires the 8 most recent states and actions to make decisions, our sampling
187 process is more complex compared to the original MLP-based PPO model. To handle this complexity,
188 we maintain a trajectory collector for each parallel environment. This collector ensures that data
189 is formatted correctly for the network from the raw trajectories. Additionally, once a trajectory is
190 completed, the collector facilitates the accumulation of data for the training process. For further
191 details, please refer to the implementation below.

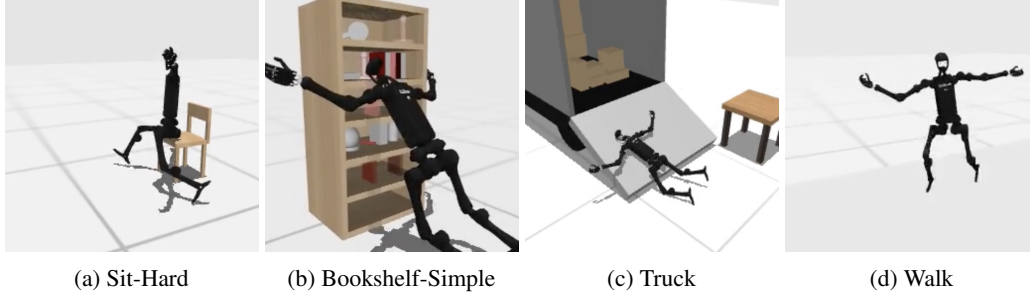


Figure 2: Some Visual Results on Four Environments

4 Experiments

To evaluate the effectiveness of our improved PPO algorithm, we conducted experiments across four distinct environments: Sit-Hard, Bookshelf-Simple, Truck, and Walk. These environments were chosen to represent a variety of tasks, ranging from basic manipulation to more complex locomotion and coordination challenges.

It is worth noting that we do not include the original PPO algorithm in our training curves, as it consistently fails to demonstrate stable progress on these tasks. This observation aligns with the findings presented in the HumanoidBench paper, where PPO’s performance was notably suboptimal, particularly on high-dimensional environments like HumanoidBench.

4.1 Sit-Hard Environment

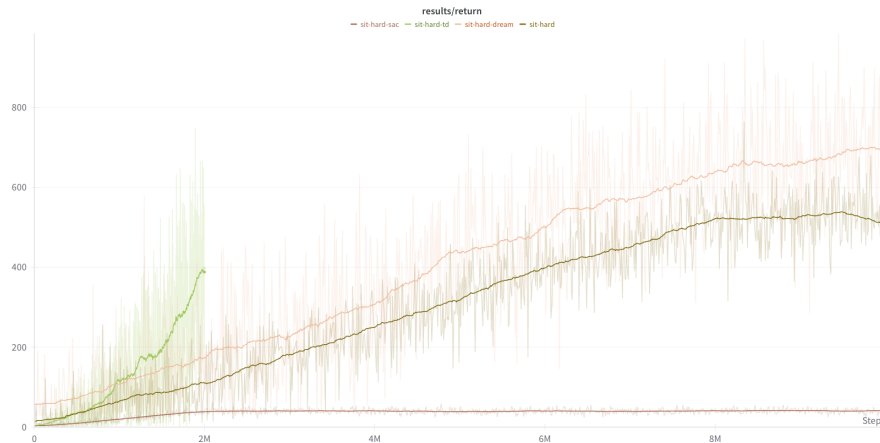


Figure 3: Reward Curves on Sit-Hard

In the Sit-Hard environment, our improved PPO model outperforms the model-free Soft Actor-Critic (SAC) algorithm by a significant margin. The performance of our algorithm reaches a level comparable to that of model-based methods such as TD-MPC2 and DreamerV3. This indicates that, despite being a model-free algorithm, our PPO variant is able to achieve near-model-based performance in this environment.

207 4.2 Bookshelf-Simple Environment

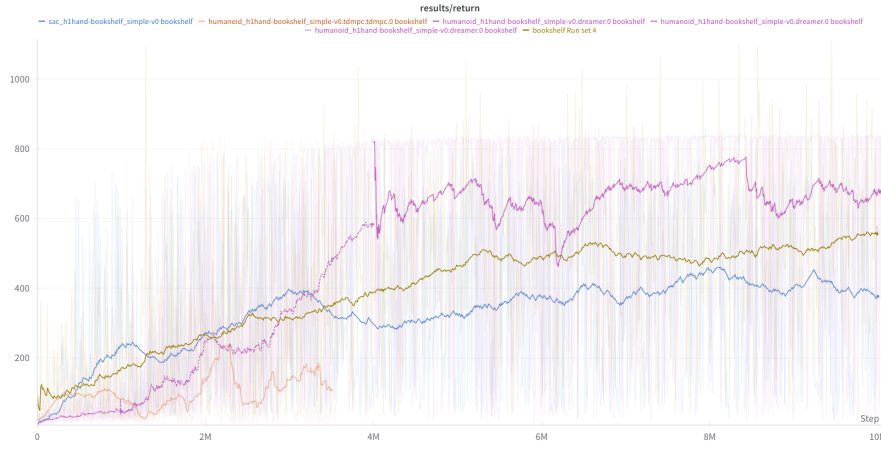


Figure 4: Reward Curves on BookShelf_Simple

208 In the Bookshelf-Simple environment, our model, SAC, and DreamerV3 exhibit similar behavior in
 209 terms of actions—they all stay close to the bookshelf, although they cannot finish the manipulation
 210 task. However, TD-MPC2 behaves differently, as it deviates from the bookshelf and strays from the
 211 intended task. This demonstrates that our PPO model, SAC, and DreamerV3 maintain a more stable
 212 and goal-directed behavior, while TD-MPC2 struggles with consistency in this environment.

213 4.3 Truck Environment

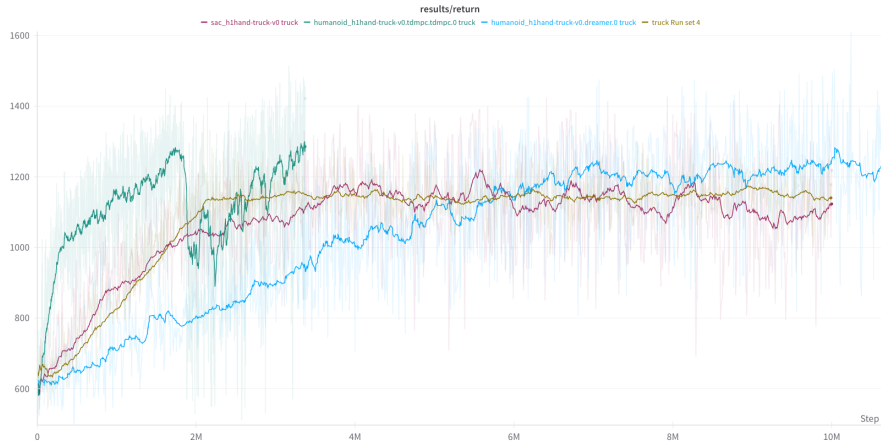


Figure 5: Reward Curves on Truck

214 The Truck environment presents a case where reward design significantly impacts learning perfor-
 215 mance. In this setting, all models—including ours—fail to achieve meaningful progress and instead
 216 simply remain stationary, lying on the ground. The rewards for this environment are poorly designed,
 217 leading to a plateau in performance for all algorithms, which highlights the importance of well-tuned
 218 reward structures in reinforcement learning tasks.

219 4.4 Walk Environment

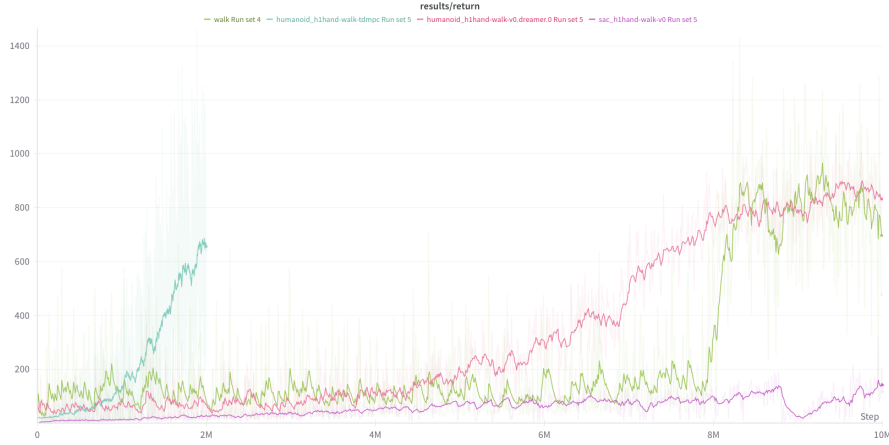


Figure 6: Reward Curves on Walk

220 In the Walk environment, SAC struggles to learn a valid walking policy, failing to make any mean-
 221 ingful progress. In contrast, both the model-based methods (TD-MPC2 and DreamerV3) and our
 222 improved PPO algorithm successfully learn to walk, but with some distinctions. Our model initially
 223 exhibits more oscillatory behavior, which can be attributed to the relatively large action entropy
 224 term we incorporate to encourage exploration during early training steps. While this results in some
 225 instability at the beginning, our PPO model quickly converges to a successful walking policy after
 226 sufficient exploration. This ability to recover and stabilize after a period of exploration demonstrates
 227 the robustness of our approach.

228 Overall, the experiments demonstrate that our method significantly improves PPO’s performance in
 229 high-dimensional, complex humanoid robot learning tasks.

230 5 Discussion And Analysis

231 6 Conclusion

232 In this paper, we present significant improvements in the data efficiency and performance of the
 233 PPO algorithm on the HumanoidBench benchmark, which involves a humanoid robot tackling
 234 complex tasks requiring both locomotion and manipulation. By enhancing the representation learning
 235 capability of PPO through sequentially modeling the policy-value network with interleaved states and
 236 actions, and incorporating techniques such as state normalization, reward scaling, and Generalized
 237 Advantage Estimation, we address the limitations highlighted in the original HumanoidBench paper.
 238 Our enhanced PPO algorithm not only surpasses traditional model-free methods, including the
 239 baseline PPO and Soft Actor-Critic algorithms, but also achieves performance levels comparable
 240 to state-of-the-art model-based approaches such as DreamerV3 and TD-MPC2. This work offers
 241 valuable insights into improving the sample efficiency and representation learning capabilities of
 242 PPO, which are critical for deploying PPO in real-world, complex environments.

243 References

- 244 [1] Petros Christodoulou. Soft actor-critic for discrete action settings. *arXiv preprint arXiv:1910.07207*, 2019.
- 245 [2] Rahul Dey and Fathi M Salem. Gate-variants of gated recurrent unit (gru) neural networks. In *2017 IEEE*
 246 *60th international midwest symposium on circuits and systems (MWSCAS)*, pages 1597–1600. IEEE, 2017.
- 247 [3] Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Firdaus Janoos, Larry Rudolph, and
 248 Aleksander Madry. Implementation matters in deep rl: A case study on ppo and trpo. In *International*
 249 *conference on learning representations*, 2019.

- 250 [4] Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Firdaus Janoos, Larry Rudolph, and
251 Aleksander Madry. Implementation matters in deep policy gradients: A case study on ppo and trpo. *arXiv*
252 *preprint arXiv:2005.12729*, 2020.
- 253 [5] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum
254 entropy deep reinforcement learning with a stochastic actor. In *International conference on machine*
255 *learning*, pages 1861–1870. PMLR, 2018.
- 256 [6] Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar,
257 Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv*
258 *preprint arXiv:1812.05905*, 2018.
- 259 [7] Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete
260 world models. *arXiv preprint arXiv:2010.02193*, 2020.
- 261 [8] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through
262 world models. *arXiv preprint arXiv:2301.04104*, 2023.
- 263 [9] Nicklas Hansen, Hao Su, and Xiaolong Wang. Td-mpc2: Scalable, robust world models for continuous
264 control. *arXiv preprint arXiv:2310.16828*, 2023.
- 265 [10] Nicklas Hansen, Xiaolong Wang, and Hao Su. Temporal difference learning for model predictive control.
266 *arXiv preprint arXiv:2203.04955*, 2022.
- 267 [11] Matthias Hutter, Kevin Mets, and Steven Latré. Hierarchical reinforcement learning: A survey
268 and open research challenges. *Machine Learning and Knowledge Extraction*, 4(1):172–221, 2022.
- 269 [12] Shubham Pateria, Budhitama Subagdja, Ah-hwee Tan, and Chai Quek. Hierarchical reinforcement learning:
270 A comprehensive survey. *ACM Computing Surveys (CSUR)*, 54(5):1–35, 2021.
- 271 [13] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional
272 continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- 273 [14] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
274 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 275 [15] Carmelo Sferrazza, Dun-Ming Huang, Xingyu Lin, Youngwoon Lee, and Pieter Abbeel. Humanoid-
276 bench: Simulated humanoid benchmark for whole-body locomotion and manipulation. *arXiv preprint*
277 *arXiv:2403.10506*, 2024.
- 278 [16] Alex Sherstinsky. Fundamentals of recurrent neural network (rnn) and long short-term memory (lstm)
279 network. *Physica D: Nonlinear Phenomena*, 404:132306, 2020.
- 280 [17] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In
281 *2012 IEEE/RSJ international conference on intelligent robots and systems*, pages 5026–5033. IEEE, 2012.