

PPO Performance in MuJoCo Environments: A Comparative Analysis

Author: BenchChem Technical Support Team. **Date:** December 2025

Compound of Interest

Compound Name: Ppo-IN-5

Cat. No.: B12371345

[Get Quote](#)

Proximal Policy Optimization (PPO) has emerged as a leading reinforcement learning algorithm for continuous control tasks, frequently benchmarked in the MuJoCo physics simulation environment.[1] Its balance of sample efficiency, stability, and ease of implementation has made it a popular choice for robotics and control applications.[2][3] This guide provides a comparative analysis of PPO's performance against other state-of-the-art algorithms in various MuJoCo tasks, supported by experimental data and detailed methodologies.

Algorithm Overview

PPO is an on-policy, policy gradient algorithm that optimizes a "surrogate" objective function using a clipped probability ratio, which restricts the size of policy updates at each training step. [3] This clipping mechanism is key to PPO's stability, preventing drastic performance drops that can occur with traditional policy gradient methods. PPO is often compared to other prominent model-free algorithms such as Soft Actor-Critic (SAC) and Twin Delayed Deep Deterministic Policy Gradient (TD3), which are off-policy methods known for their sample efficiency.[4]

Performance Benchmark

The following tables summarize the performance of PPO and other leading reinforcement learning algorithms on a selection of MuJoCo environments. The performance is typically measured as the average return over a number of episodes after a fixed number of training timesteps.

Note: The results presented below are aggregated from various studies and benchmark reports. Direct comparison can be challenging due to slight variations in experimental setups.

Environment	PPO	SAC	TD3
HalfCheetah-v2	~4000 - 8000	~10000 - 12000	~9000 - 11000
Hopper-v2	~2500 - 3500	~3500 - 3800	~3400 - 3700
Walker2d-v2	~3000 - 5000	~4500 - 5500	~4000 - 5000
Ant-v2	~2500 - 4500	~5000 - 6000	~4500 - 5500

Table 1: Comparative performance of PPO, SAC, and TD3 on select MuJoCo environments. Values represent the approximate range of average returns. Higher values indicate better performance. Bolded values indicate the generally top-performing algorithm for that environment.

Environment	PPO Average Return	PPO Standard Deviation
HalfCheetah-v2	7534	1354
Hopper-v2	3478	213
Walker2d-v2	4891	572
Ant-v2	3982	1123

Table 2: Example PPO performance from a specific benchmark run, showcasing average return and standard deviation after 3 million timesteps. These values can vary based on implementation and hyperparameter tuning.

Experimental Protocols

Reproducibility of reinforcement learning experiments is crucial. Below are typical experimental setups used for benchmarking PPO and other algorithms in MuJoCo environments.

Common Hyperparameters for PPO:

Hyperparameter	Value	Description
Discount Factor (γ)	0.99	The factor by which future rewards are discounted.
GAE Parameter (λ)	0.95	The parameter for Generalized Advantage Estimation, balancing bias and variance.
Clipping Parameter (ϵ)	0.2	The clipping range for the surrogate objective function.
Epochs per Update	10	The number of epochs of stochastic gradient ascent to perform on the collected data.
Minibatch Size	64	The size of minibatches for the stochastic gradient ascent updates.
Optimizer	Adam	The optimization algorithm used.
Learning Rate	3e-4 (often annealed)	The learning rate for the optimizer.
Value Function Coef.	0.5	The weight of the value function loss in the total loss.
Entropy Coefficient	0.0	The weight of the entropy bonus, encouraging exploration.

Network Architecture:

The policy and value functions are commonly represented by feedforward neural networks. A typical architecture for MuJoCo tasks consists of:

- Two hidden layers with 64 units each.
- Tanh activation functions for the hidden layers.

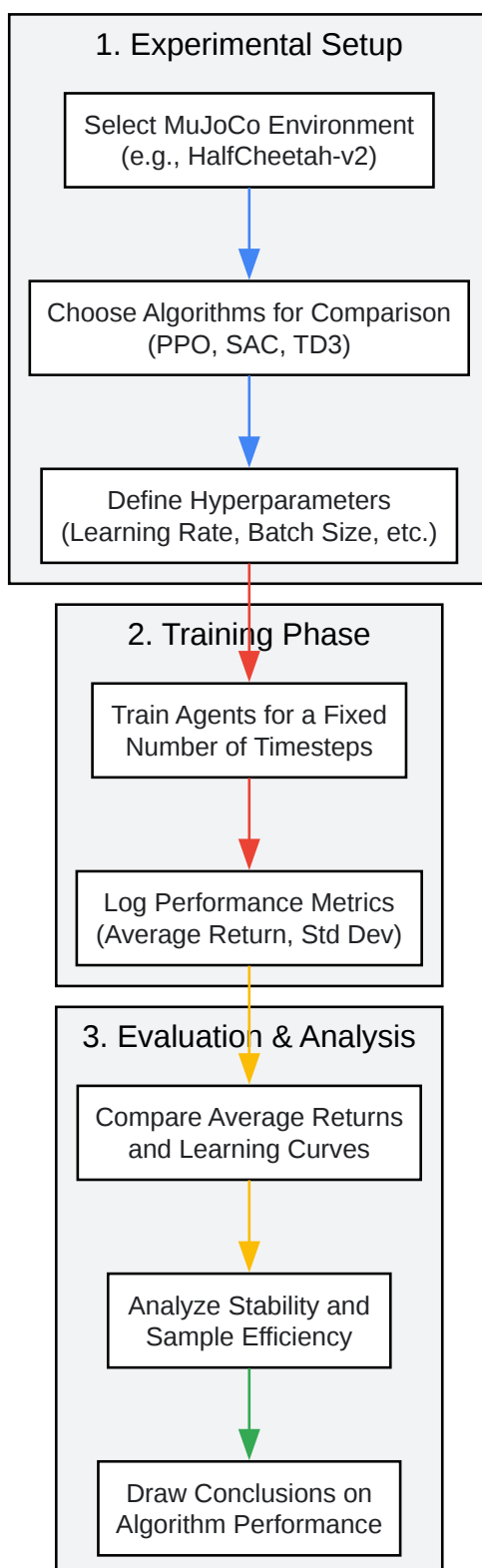
- The policy network outputs the mean of a Gaussian distribution for each action dimension, with state-independent standard deviations that are also learned.

Data Collection:

- On-policy data collection: PPO collects a batch of experience by running the current policy in the environment.
- Number of steps: A common practice is to collect 2048 or 4096 steps of agent-environment interaction per update.
- Normalization: Observations and advantages are often normalized to improve training stability.

Logical Workflow for PPO Performance Evaluation

The following diagram illustrates the typical workflow for benchmarking PPO performance in a MuJoCo environment.



[Click to download full resolution via product page](#)

PPO Benchmarking Workflow

Conclusion

PPO consistently demonstrates strong and stable performance across a variety of MuJoCo continuous control tasks. While off-policy algorithms like SAC and TD3 may achieve higher final returns in some environments due to their improved sample efficiency, PPO remains a robust and reliable baseline. Its relative simplicity and stability make it an excellent choice for a wide range of research and development applications. For researchers and professionals, the choice between PPO and its off-policy counterparts will often depend on the specific requirements of the task, including the importance of sample efficiency versus training stability and ease of implementation.

Need Custom Synthesis?

BenchChem offers custom synthesis for rare earth carbides and specific isotopic labeling.

Email: info@benchchem.com or [Request Quote Online](#).

References

- 1. medium.com [medium.com]
- 2. GitHub - danimatasd/MUJOCO-AIDL: Reinforced learning on Mujoco for AIDL final project [github.com]
- 3. Proximal Policy Optimization — Spinning Up documentation [spinningup.openai.com]
- 4. Sim-to-Real: A Performance Comparison of PPO, TD3, and SAC Reinforcement Learning Algorithms for Quadruped Walking Gait Generation [scirp.org]
- To cite this document: BenchChem. [PPO Performance in MuJoCo Environments: A Comparative Analysis]. BenchChem, [2025]. [Online PDF]. Available at: [\[https://www.benchchem.com/product/b12371345#ppo-performance-benchmark-on-mujoco-environments\]](https://www.benchchem.com/product/b12371345#ppo-performance-benchmark-on-mujoco-environments)

Disclaimer & Data Validity:

The information provided in this document is for Research Use Only (RUO) and is strictly not intended for diagnostic or therapeutic procedures. While BenchChem strives to provide

accurate protocols, we make no warranties, express or implied, regarding the fitness of this product for every specific experimental setup.

Technical Support: The protocols provided are for reference purposes. Unsure if this reagent suits your experiment? [[Contact our Ph.D. Support Team for a compatibility check](#)]

Need Industrial/Bulk Grade? [Request Custom Synthesis Quote](#)

BenchChem

Our mission is to be the trusted global source of essential and advanced chemicals, empowering scientists and researchers to drive progress in science and industry.

Contact

Address: 3281 E Guasti Rd
Ontario, CA 91761, United States
Phone: (601) 213-4426
Email: info@benchchem.com