

# A Comparative Analysis of PPO-Clip and PPO-Penalty in Reinforcement Learning

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

## Compound of Interest

Compound Name: Ppo-IN-5

Cat. No.: B12371345

[Get Quote](#)

In the landscape of reinforcement learning, Proximal Policy Optimization (PPO) has emerged as a robust and widely adopted algorithm for policy optimization. Its appeal lies in its blend of sample efficiency, stability, and ease of implementation. PPO navigates the crucial trade-off between taking sufficiently large policy update steps to ensure learning progress and avoiding overly aggressive updates that can lead to performance collapse. This is achieved through two primary variants: PPO-Clip and PPO-Penalty. This guide provides a comprehensive comparison of these two methods, supported by experimental data and detailed protocols, to aid researchers and practitioners in selecting the appropriate variant for their needs.

## Core Concepts: PPO-Clip vs. PPO-Penalty

Both PPO-Clip and PPO-Penalty strive to keep the new policy close to the old one, but they employ different mechanisms to achieve this goal.

PPO-Clip, the more prevalent variant, utilizes a clipped surrogate objective function.<sup>[1]</sup> This function constrains the policy update by clipping the probability ratio between the new and old policies.<sup>[1]</sup> This simple yet effective mechanism prevents the new policy from deviating too far from the previous one, thereby enhancing training stability.<sup>[2]</sup> Its popularity stems from its straightforward implementation and strong empirical performance across a variety of tasks.<sup>[1]</sup>

PPO-Penalty, on the other hand, incorporates a soft constraint on the policy update by adding a penalty term to the objective function. This penalty is proportional to the Kullback-Leibler (KL) divergence between the new and old policies.<sup>[1]</sup> An adaptive coefficient for this penalty term is

typically used, which is adjusted based on the observed KL divergence during training. This allows for more explicit control over the magnitude of policy changes.

## Performance Benchmark

The following table summarizes the performance of PPO-Clip and an adaptive KL penalty approach (akin to PPO-Penalty) on several continuous control benchmarks from the MuJoCo suite, as presented in the original Proximal Policy Optimization paper. The scores are normalized, where 0 corresponds to the performance of a random policy and 1 corresponds to the performance of Trust Region Policy Optimization (TRPO).

Environment	PPO-Clip ( $\epsilon=0.2$ )	PPO with Adaptive KL Penalty
Average	0.83	0.70
HalfCheetah	0.77	0.65
Hopper	0.90	0.82
InvertedDoublePendulum	0.75	0.60
InvertedPendulum	0.95	0.92
Reacher	0.65	0.55
Swimmer	0.92	0.88
Walker2d	0.85	0.75

Note: The results are based on the findings reported in the original PPO paper by Schulman et al. (2017). The values represent the average normalized scores over 21 runs of the algorithm on 7 environments.

The empirical results suggest that PPO-Clip generally outperforms the adaptive KL penalty variant across a range of continuous control tasks. The simplicity and effectiveness of the

clipping mechanism often lead to more stable and higher-performing policies.

## Experimental Protocols

Reproducing benchmark results requires a clear understanding of the experimental setup. The following protocols are based on common practices for benchmarking PPO variants in continuous control environments.

### PPO-Clip

- Objective Function: Clipped Surrogate Objective.
- Hyperparameters:
  - Clipping parameter ( $\epsilon$ ): Typically set to 0.2. This parameter defines the range  $[1-\epsilon, 1+\epsilon]$  within which the probability ratio is clipped.
  - Discount factor ( $\gamma$ ): Commonly set to 0.99.
  - GAE parameter ( $\lambda$ ): Usually set to 0.95 for Generalized Advantage Estimation.
  - Number of epochs: Typically between 3 and 15.
  - Minibatch size: A common choice is 64.
  - Learning rate: Often in the range of  $3e-4$ , potentially with linear decay.
  - Entropy coefficient: A small value, such as 0.01, is often used to encourage exploration.
- Network Architecture: For MuJoCo tasks, a common architecture consists of two hidden layers with 64 units each and tanh activation functions. The policy and value functions may or may not share parameters.
- Optimization: The Adam optimizer is typically used.

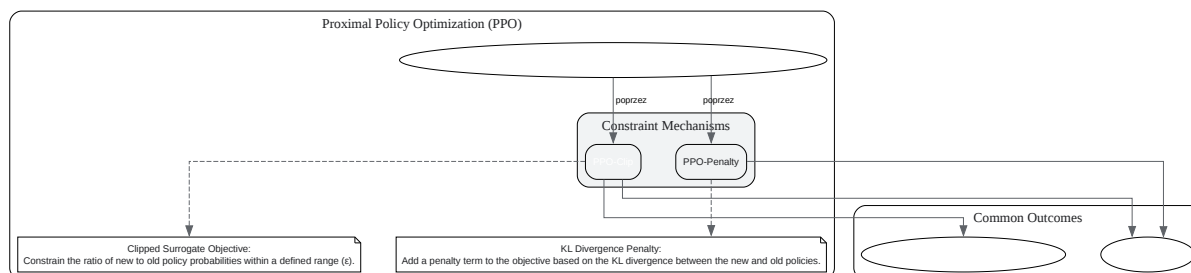
### PPO-Penalty

- Objective Function: Surrogate Objective with a KL Penalty term.

- Hyperparameters:
  - Initial KL penalty coefficient ( $\beta$ ): A starting value, often around 1.0, is chosen.
  - Target KL divergence: A target value for the KL divergence between the old and new policies is set, for example, 0.01. The penalty coefficient is then adapted based on whether the observed KL divergence is higher or lower than this target.
  - Other hyperparameters such as the discount factor, GAE parameter, number of epochs, minibatch size, and learning rate are typically in the same range as for PPO-Clip.
- Network Architecture: Similar to PPO-Clip, a feedforward neural network with two hidden layers of 64 units and tanh activations is a common choice for both the policy and value functions.
- Optimization: Adam is the standard optimizer.

## Logical Relationship and Algorithmic Flow

The fundamental difference between PPO-Clip and PPO-Penalty lies in how they constrain the policy update. The following diagram illustrates this logical relationship.



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

### PPO Variants: Core Mechanisms

## Conclusion

Both PPO-Clip and PPO-Penalty are effective methods for stable and efficient policy optimization in reinforcement learning. PPO-Clip is generally favored for its simplicity, ease of tuning, and strong empirical performance, making it a go-to algorithm for a wide range of applications. PPO-Penalty offers a more explicit way to control policy divergence through the KL penalty, which can be beneficial in scenarios where precise control over the policy update is critical. The choice between the two variants will depend on the specific requirements of the task, including the desired level of implementation complexity and the need for explicit control over policy updates. For most practical applications, PPO-Clip provides a robust and high-performing solution.

### *Need Custom Synthesis?*

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

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

## References

- 1. aiarts.medium.com [aiarts.medium.com]
- 2. Proximal Policy Optimization — Spinning Up documentation [spinningup.openai.com]
- To cite this document: BenchChem. [A Comparative Analysis of PPO-Clip and PPO-Penalty in Reinforcement Learning]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b12371345#benchmarking-ppo-variants-like-ppo-clip-and-ppo-penalty]

---

### 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](mailto:info@benchchem.com)