

# Generated Model Performance Report

Simon Green

February 7, 2025

## 1 Results and Model Comparison

This report presents the performance evaluation of four reinforcement learning models: TRPO, PPO, TRPOER, and TRPOR. These models are implemented using **Stable Baselines3** and utilize **mini-batch gradient descent** for optimization, ensuring efficient and stable updates during training. The entropy calculations guiding the models are based on the entropy of each batch, which influences regularization and experience replay mechanisms.

### 1.1 TRPO (Trust Region Policy Optimization)

Originally proposed by Schulman et al. [2], TRPO is a policy gradient method that constrains updates using a trust region to ensure stability in training.

### 1.2 PPO (Proximal Policy Optimization)

Introduced by Schulman et al. [3], PPO improves upon TRPO by using a clipped surrogate objective to ensure efficient and stable policy updates.

### 1.3 TRPOR (TRPO with Entropy Regularization)

This model extends TRPO by introducing entropy regularization only in the policy objective. The entropy coefficient hyperparameter guides the degree of regularization, ensuring a balance between exploration and exploitation. The entropy guiding this model is computed at the batch level, dynamically adjusting policy updates.

### 1.4 TRPOER (TRPO with Entropy Regularized Experience Replay)

This model extends TRPO by incorporating entropy-based experience replay and an additional policy entropy regularization term. It utilizes a prioritized experience replay buffer sampled according to batch entropy values and a hyperparameter coefficient. The method enables bidirectional adaptive sampling, adjusting both the number and direction of sampled experiences to optimize learning. The adaptive sampling function is formulated as:

$$S = \text{clip} \left( (M - m) \times \begin{cases} 1 - \left| \frac{H}{|\lambda + \epsilon|} \right|, & \lambda > 0 \\ \left| \frac{H}{|\lambda + \epsilon|} \right|, & \lambda < 0 \end{cases} + m, m, M \right) \quad (1)$$

where  $S$  is the number of samples,  $H$  represents batch entropy,  $\lambda$  is the sampling coefficient, and  $M, m$  are the maximum and minimum sample limits.

## 1.5 GenTRPO (Generative Experience Replay Trust Region Policy Optimization with Entropy Regularization)

Quite a mouth full, we’ll find a better name. The GenTRPO algorithm extends the Trust Region Policy Optimization with Entropy Regularization (TRPOER) framework [2, 3] by incorporating a generative model to augment the experience replay buffer. The key idea is to leverage synthetic experiences generated by a forward dynamics model to complement real experiences, thereby improving exploration and sample efficiency.

In the GenTRPO framework, the experiences used for policy updates are sampled from a replay buffer. The sampling strategy ensures that half of the samples in each batch are real experiences collected from the environment, while the other half are generated by the forward dynamics model. This combination of real and synthetic data balances model fidelity with exploratory richness, enabling the policy to generalize effectively while maintaining stability during optimization.

The generative component of GenTRPO relies on a forward dynamics model inspired by the intrinsic curiosity module [1]. The forward dynamics model comprises an encoder and a dynamics predictor. The encoder maps raw states  $s$  into a compact latent space representation  $h(s)$ , capturing the essential features of the environment. The dynamics predictor then takes the latent state  $h(s)$  and action  $a$  as input and predicts the next latent state  $h(s')$ , effectively modeling the transition function  $P(s'|s, a)$ . The error of this model, expressed as

$$\mathcal{F}(s, a, s', r) = \frac{1}{2} \|g(h(s), a) - h(s')\|^2 \quad (2)$$

where  $g(h(s), a)$  is the predicted latent state,  $h(s')$  is the true latent state, and  $\|\cdot\|$  represents the Euclidean norm. This error quantifies how accurately the forward dynamics model predicts the latent state transitions. It is used to compute intrinsic motivation, encouraging the agent to explore transitions that are harder to predict, thereby fostering exploration [4].

## 2 Model Performance Table

The table below summarizes the models’ performance in terms of mean and standard deviation of rewards, along with maximum and minimum rewards recorded during training. A higher mean reward suggests better overall performance, while lower standard deviation indicates increased stability.

	Ant-v5	Pendulum-v1	InvertedDouble Pendulum-v5	Humanoid-v5
PPO	1173.61M 515.51 $\mu \pm 326.48\sigma$ 1895E, 10R	−0.02M −207.56 $\mu \pm 118.01\sigma$ 2892E, 10R	9359.93M 1090.87 $\mu \pm 2910.66\sigma$ 6965E, 10R	1245.44M 464.88 $\mu \pm 298.07\sigma$ 15385E, 9R
TRPO	1960.61M 1328.68 $\mu \pm 284.88\sigma$ 1067E, 10R	−0.20M −141.97 $\mu \pm 115.55\sigma$ 2600E, 10R	9359.78M 7516.76 $\mu \pm 3879.62\sigma$ 5767E, 10R	1247.77M 369.87 $\mu \pm 311.76\sigma$ 9569E, 10R
TRPOER	1349.85M 639.65 $\mu \pm 465.91\sigma$ 1882E, 8R	−0.07M −173.49 $\mu \pm 156.48\sigma$ 2859E, 10R	9351.73M 2673.08 $\mu \pm 2620.77\sigma$ 5070E, 9R	897.72M 371.33 $\mu \pm 359.99\sigma$ 2799E, 4R
TRPOR	3590.62M 1662.87 $\mu \pm 1171.26\sigma$ 1655E, 7R	−0.22M −197.32 $\mu \pm 215.34\sigma$ 4508E, 10R	9357.32M 8045.50 $\mu \pm 2918.29\sigma$ 9123E, 5R	1416.67M 620.72 $\mu \pm 381.70\sigma$ 9225E, 10R

## 3 Performance Analysis Through Plots

The following plots visualize different aspects of model performance.

## Resampled Rewards and Outlier Removal

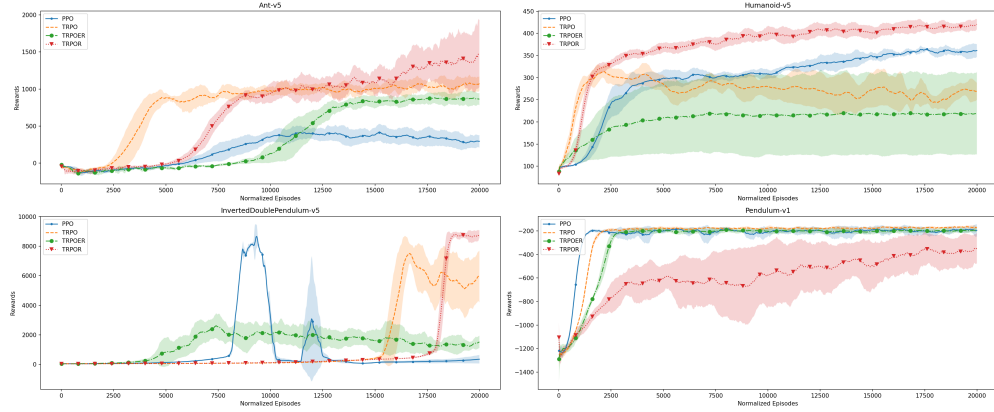


Figure 1: Resampled Rewards with Outlier Removal

This plot presents reward distributions after applying smoothing and outlier removal techniques, filtering out misleading fluctuations.

## Learning Stability

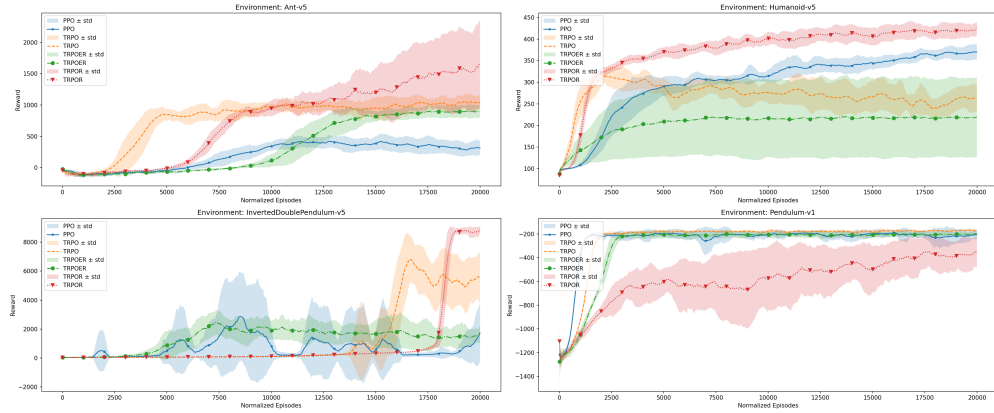


Figure 2: Learning Stability for Different Models

Learning stability is evaluated based on the smoothness of the reward curve. A more stable learning process exhibits a steadily increasing reward trajectory, whereas high variance suggests instability due to sensitivity to hyperparameters.

## Learning Stability (Coefficient of Variation)

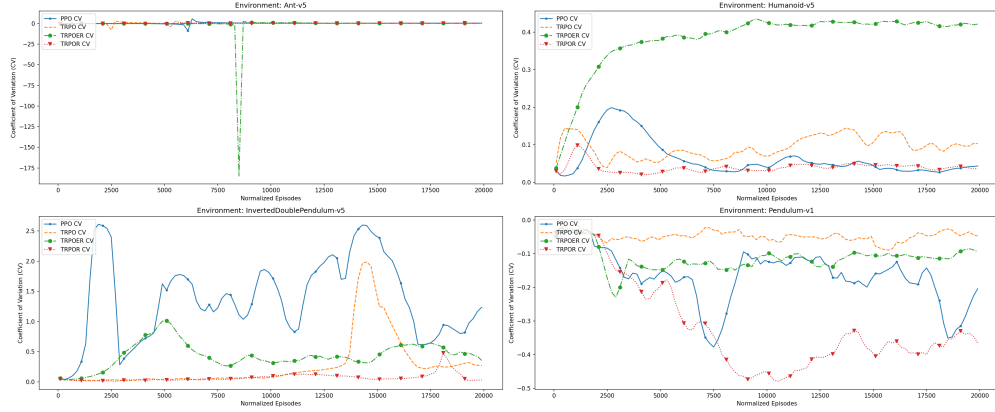


Figure 3: Learning Stability (Coefficient of Variation)

The coefficient of variation (CV) provides a normalized measure of stability. A lower CV signifies less volatile performance, whereas a higher CV indicates inconsistency due to randomness in training.

## Sample Efficiency

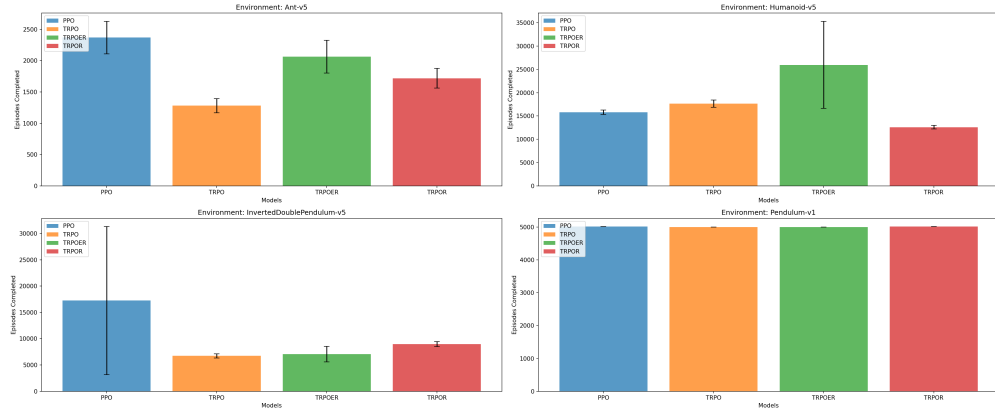


Figure 4: Sample Efficiency Across Models

Sample efficiency measures how quickly a model improves with limited training episodes. Higher sample efficiency is desirable, especially in data-scarce scenarios.

## Combined Sample Efficiency Results

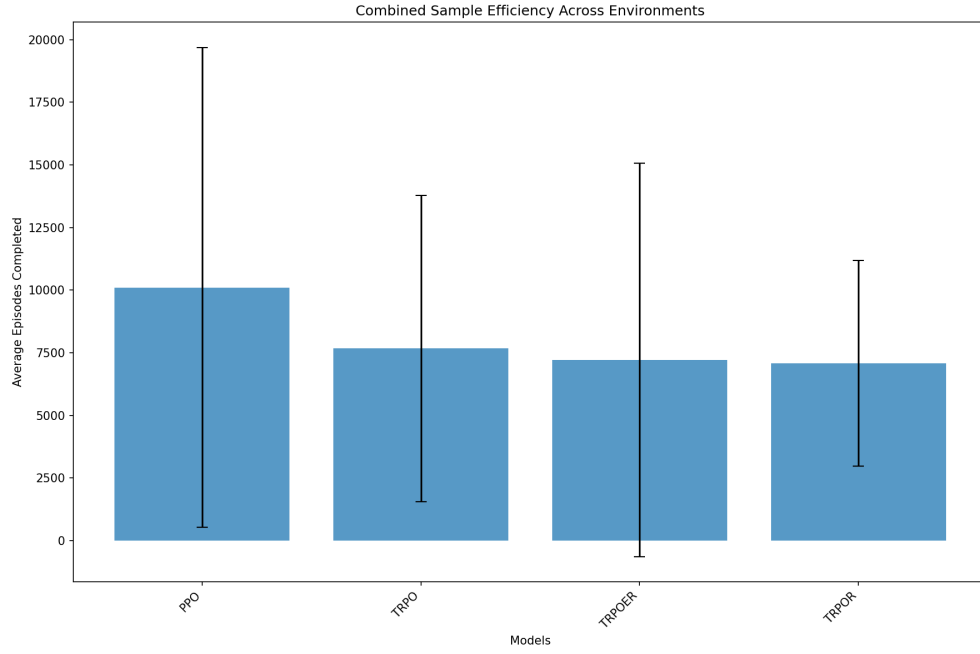


Figure 5: Combined Sample Efficiency Results

The combined sample efficiency plot aggregates results across all environments, showing how different models perform in terms of data efficiency.

## Raw Data

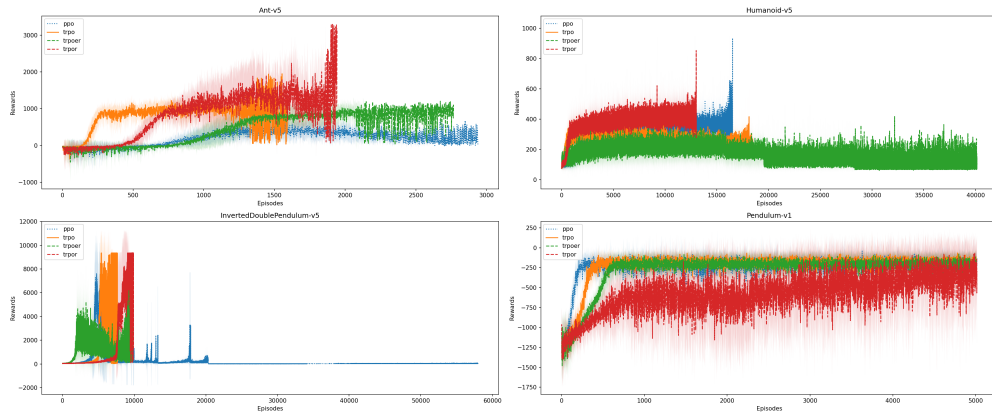


Figure 6: Raw Reward Data for Different Models

The raw data plot displays the recorded reward values without any smoothing. It provides insights into the actual training process and variability in rewards.

## References

- [1] Deepak Pathak, Pulkit Agrawal, Alexei A. Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction, 2017.
- [2] John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, and Pieter Abbeel. Trust region policy optimization, 2017.
- [3] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017.
- [4] Renhao Wang, Kevin Frans, Pieter Abbeel, Sergey Levine, and Alexei A. Efros. Prioritized generative replay, 2024.