

Group Based Reward Protection for Detecting and Mitigating Reward Hacking

Réka Erdész^{a,b}, Ede Troll^a

^aEszterházy Károly Catholic University
erdesz.reka@uni-eszterhazy.hu
troll.ede@uni-eszterhazy.hu

^bEKCU Data Science and Intelligent Systems Research Group

Abstract

In Reinforcement Learning (RL) [12], reward hacking remains a persistent challenge where agents exploit specification loopholes [2, 6, 8] rather than maximizing true objectives. While Proximal Policy Optimization (PPO) [9] is widely used, its sensitivity to random seeds often masks exploitative strategies [3], creating significant hurdles for AI safety [1, 7]. Although methods like invariant testing exist [4, 5, 11], they typically require costly post-hoc interventions.

DeepSeek’s Group Relative Policy Optimization (GRPO) [10] addresses this in Large Language Models by dispensing with the critic model and optimizing via group dynamics. This approach aggregates behavior across response groups, significantly reducing variance and offering a robust baseline for stable performance.

Inspired by GRPO, our Group Based Reward Protection (GBRP) algorithm inverts the group-relative principle: rather than optimizing for reasoning consensus, it utilizes group-level statistics to detect behavioral anomalies. This approach adapts the robustness of group dynamics to identify specification gaming without requiring external supervision.

This presentation proposes GBRP, an environment-agnostic group-based algorithm that filters reward signals in real-time. Unlike prior methods requiring ground-truth supervision, GBRP relies solely on statistical derivatives of the learning process itself (reward gradients, value loss variance). We validate our method on various standard OpenAI Reinforcement Learning environments known for reward

hacking, demonstrating that GBRP prevents “looping” strategies where standard PPO fails. In this abstract, we demonstrate one such environment, published in OpenAI’s seminal work on reward hacking [2]. The framework operates exclusively on raw rewards, aggregates behavior at the policy-group level, and employs a multi-metric detection mechanism combined with an adaptive forward-looking penalty to stabilize learning dynamics. Unlike methods requiring supervision, GBRP relies solely on statistical derivatives of the learning process (e.g., value loss variance). The framework aggregates metrics from rollout batches into groups (N) and compares them against a rolling baseline (W) using a hierarchical voting system. Level 1 monitors optimization instability (Reward Spikes, Entropy, Value loss Proxy), while Level 2 detects semantic gaming, such as Reward-Length Coupling (stalling) and Action Imbalance (repetitive actuation). Upon detecting anomalies, GBRP applies an adaptive penalty based on the metrics ($\lambda < 1.0$) to stabilize learning. The λ value is based on detection.

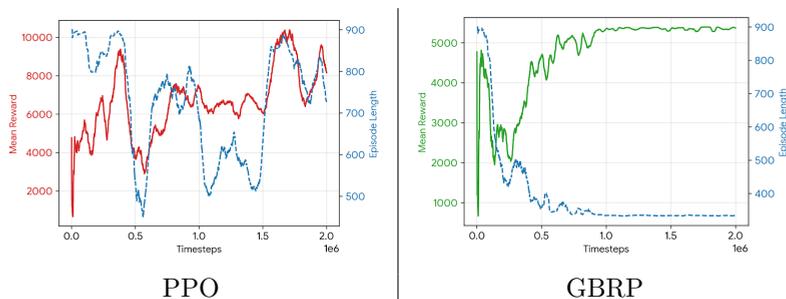


Figure 1. Agent seed = 42, 2 million max_steps, $N = 5$, $W = 3$

Figure 1 shows the metrics of the Agent’s performance. The PPO diagram (left) visually confirms the hacking strategy. It shows a synchronized rise in reward and time, indicating it ‘farms’ points by staying alive longer. The GBRP agent (right) breaks this link: reward goes up, but time stays down.

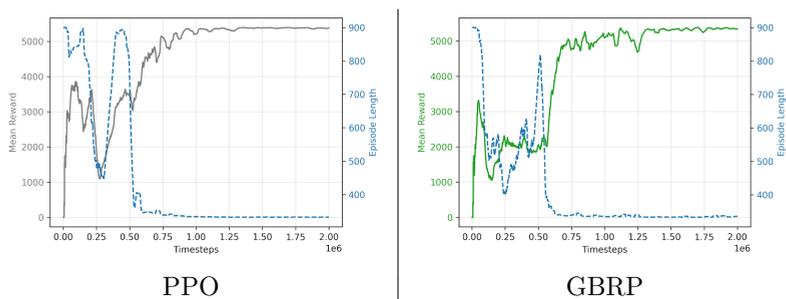


Figure 2. Agent seed 101010, 2 million max_steps, $N = 5$, $W = 3$

The control experiment in Figure 2 (seed 101010) confirms that in the absence

of anomalies, the learning dynamics of GBRP are the same as PPO. The security layer remains transparent and thus does not cause performance degradation or unnecessary overhead under normal operation.

In summary, these comparative results demonstrate GBRP’s dual capability: it effectively mitigates specification gaming (as seen in seed 42) via statistical anomaly detection, while maintaining complete transparency in benign settings (seed 101). By leveraging group dynamics rather than external supervision, the framework offers a scalable, model-agnostic path toward robust AI safety.

References

- [1] D. AMODEI, C. OLAH, J. STEINHARDT, P. CHRISTIANO, J. SCHULMAN, D. MANÉ: *Concrete Problems in AI Safety* (June 2016), DOI: [10.48550/arXiv.1606.06565](https://doi.org/10.48550/arXiv.1606.06565).
- [2] J. CLARK, D. AMODEI: *Faulty reward functions in the wild*, Accessed: 2026-01-16, OpenAI, Dec. 2016, URL: <https://openai.com/index/faulty-reward-functions/>.
- [3] L. ENGSTROM, A. ILYAS, S. SANTURKAR, D. TSIPRAS, F. JANOOS, L. RUDOLPH, A. MADRY: *Implementation Matters in Deep Policy Gradients: A Case Study on PPO and TRPO*, in: International Conference on Learning Representations, 2020.
- [4] T. EVERITT, M. HUTTER, R. KUMAR, V. KRAKOVNA: *Reward tampering problems and solutions in reinforcement learning: A causal influence diagram perspective*, Synthese 198.27 (2021), pp. 6429–6461, DOI: <https://doi.org/10.1007/s11229-021-03141-4>.
- [5] J. GARCÍA, F. FERNÁNDEZ: *A Comprehensive Survey on Safe Reinforcement Learning*, Journal of Machine Learning Research 16.42 (2015), pp. 1437–1480, URL: <http://jmlr.org/papers/v16/garcia15a.html>.
- [6] V. KRAKOVNA, J. UESATO, V. MIKULIK, M. RAHTZ, T. EVERITT, R. KUMAR, Z. KENTON, J. LEIKE, S. LEGG: *Specification gaming: the flip side of AI ingenuity*, Accessed: 2026-01-16, Google DeepMind, Apr. 2020, URL: <https://deepmind.google/blog/specification-gaming-the-flip-side-of-ai-ingenuity/>.
- [7] J. LEIKE, M. MARTIC, V. KRAKOVNA, P. A. ORTEGA, T. EVERITT, A. LEFRANCQ, L. ORSEAU, S. LEGG: *AI Safety Gridworlds*, 2017, arXiv: [1711.09883](https://arxiv.org/abs/1711.09883) [cs.LG], URL: <https://arxiv.org/abs/1711.09883>.
- [8] A. PAN, K. BHATIA, J. STEINHARDT: *The Effects of Reward Misspecification: Mapping and Mitigating Misaligned Models*, in: International Conference on Learning Representations (ICLR), 2022.
- [9] J. SCHULMAN, F. WOLSKI, P. DHARIWAL, A. RADFORD, O. KLIMOV: *Proximal Policy Optimization Algorithms* (2017), arXiv: [1707.06347](https://arxiv.org/abs/1707.06347) [cs.LG], URL: <https://arxiv.org/abs/1707.06347>.
- [10] Z. SHAO, P. WANG, Q. ZHU, R. XU, J. SONG, X. BI, H. ZHANG, M. ZHANG, Y. K. LI, Y. WU, D. GUO: *DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models*, 2024, arXiv: [2402.03300](https://arxiv.org/abs/2402.03300) [cs.CL], URL: <https://arxiv.org/abs/2402.03300>.
- [11] I. F. SHIHAB, S. AKTER, A. SHARMA: *Detecting and Mitigating Reward Hacking in Reinforcement Learning Systems: A Comprehensive Empirical Study*, arXiv preprint arXiv:2507.05619 (July 2025), Accessed: 2026-01-16, URL: <https://arxiv.org/pdf/2507.05619>.
- [12] R. S. SUTTON, A. G. BARTO: *Reinforcement Learning: An Introduction*, 2nd, MIT Press, 2018, URL: <https://web.stanford.edu/class/psych209/Readings/SuttonBartoIPRLBook2ndEd.pdf>.