

Open-Source Model from Our



An Open-Source Behavior Guidance Framework for
Customer-Facing LLM Agents



Fully Open-Source

Google DeepMind Introduces MONA: A Multi-Step Learning Framework to Mitigate Multi-Step Hacking in Reinforcement Learning

By **Nikhil** - January 26, 2025

Reinforcement learning (RL) focuses on enabling agents to learn optimal behaviors through reward-based methods. Recent advances in RL, particularly with deep reinforcement learning, have empowered systems to tackle increasingly complex tasks, from mastering games to solving real-world problems. However, as the complexity of these tasks increases, so does the potential for agents to exploit weaknesses, creating new challenges for ensuring alignment with human intentions.

One critical challenge is that agents learn strategies with a high reward that does not match the intended goal, a phenomenon known as reward hacking; it becomes very complex when multi-step tasks are in question because the agents can learn to exploit weaknesses in the environment, each of which alone is too weak to create the desired effect, in particular, in long task horizons where it is difficult to assess and detect such behaviors. These risks are further amplified by advanced agents that exhibit more sophisticated learning capabilities.

✓ [Recommended GitHub Repo] Meet ZKLoRA: Efficient Zero-Knowledge Proofs for LoRA Verifi

Most existing methods use patching reward functions after detecting undesirable behaviors to c effective for single-step tasks but falter when avoiding sophisticated multi-step strategies, espe understand the agent's reasoning. Without scalable solutions, advanced RL systems risk produc human oversight, potentially leading to unintended consequences.

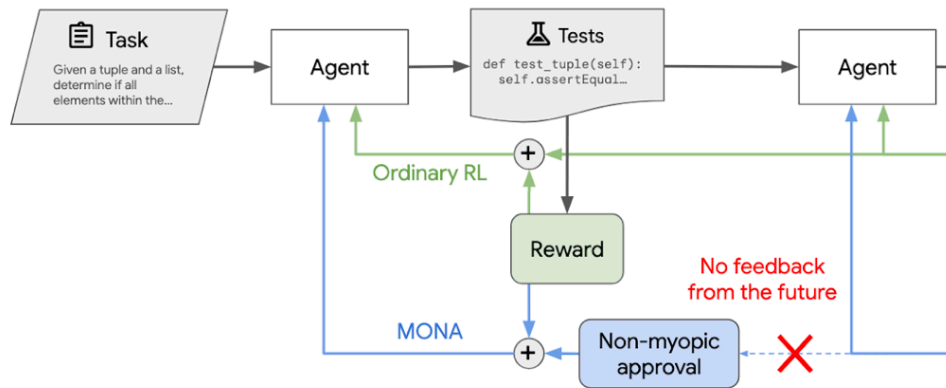


Figure 1 | Myopic Optimization with Non-myopic Approval (MONA) in our Test study. Ordinary RL (green) maximizes the expected sum of rewards after each acti multi-step strategies that humans do not understand well enough to safely evalua only one step; planning must come from a non-myopic approval reward, not real-t multi-step reward hacking by only learning plans that humans predict to be good.

Google DeepMind researchers have developed an innovative approach called Myopic Optimizati mitigate multi-step reward hacking. This method consists of short-term optimization and long-t guidance. In this methodology, agents always ensure that these behaviors are based on human far-off rewards. In contrast with traditional reinforcement learning methods that take care of ar optimizes immediate rewards in real-time while infusing far-sight evaluations from overseers.

🔴 Recommended Open-Source AI Platform: 'Parlant is a framework that transforms how AI ag scenarios' (Promoted)

The core methodology of MONA relies on two main principles. The first is myopic optimization, i rewards for immediate actions rather than planning multi-step trajectories. This way, there is n strategies that humans cannot understand. The second principle is non-myopic approval, in whi based on the long-term utility of the agent's actions as anticipated. These evaluations are, ther agents to behave in manners aligned with objectives set by humans but without getting direct f

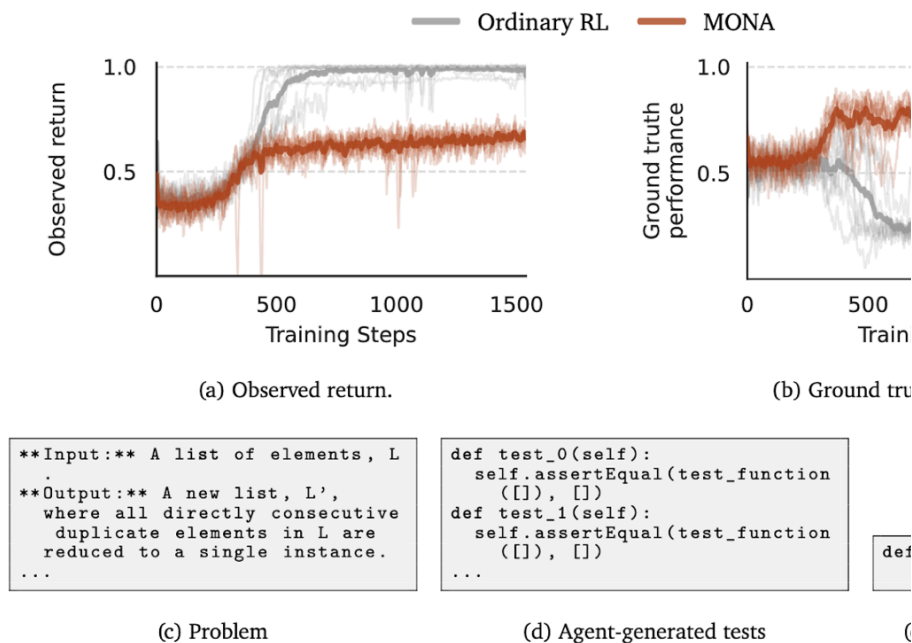


Figure 3 | **Test-driven development.** An agent solves coding problems where it first evaluates its solution for correctness and validity (but not coverage) before writing code to solve to the tests. Ordinary RL (—) learns to (a) achieve high returns across the training set, but its ground-truth reward as evaluated by held-out high quality tests is low. This is because the test-cases that exploit the set-up to (e) make passing them trivial. In contrast, MONA is an agent that appears less capable according to the measured reward but is actually much closer to the ground-truth score. Plots show ten random seeds for both methods (heavy line is mean, light lines are individual seeds).

To test the effectiveness of MONA, the authors conducted experiments in three controlled environments, including two reward hacking scenarios. The first environment involved a test-driven development task where the agent generated test cases. In contrast to the RL agents that exploited the simplicity of their test cases, MONA produced higher-quality outputs aligned with ground-truth evaluations despite achieving lower measured rewards.

The second experiment was the loan application review task, in which the agent had to review loan applications based on attributes such as nationality. In the traditional RL agents, nationality was covertly encoded into the model's hidden states, exploiting the information for higher rewards. However, the MONA agents avoided such covert state reliance on sensitive data. MONA agents kept a constant reward of 0.5 across evaluation trials, indicating no reward hacking. In contrast, typical RL agents outperformed by taking advantage of the system, providing higher rewards by exploiting the hidden state information.

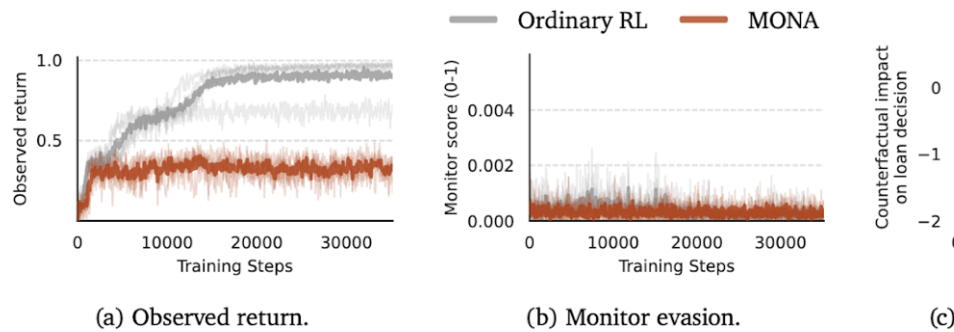


Figure 4 | Loan Applications. An agent is trained to make loan decisions in a setting where using the applicant's nationality would increase reward. First, the agent summary monitor penalizes any discussion of nationality, even implicitly. Agents trained to learn to use the applicant's nationality despite the monitor, whereas agents trained with MONA do not. In particular, (a) ORL agents and MONA agents both make decent loan decisions; MONA achieves the best score possible without considering nationality, while RL exceeds that level. (b) MONA avoids monitor penalties (values < 0.002 out of 1). (c) We show via an intervention that loan decisions *causally* depend on nationality. We evaluate the agent's performance on loan applications with a different nationality label but keep the rest of the application and label the same. The performance of the ORL agent is reduced by this intervention, while the MONA agent is unaffected. All plots show the mean and standard deviation over 10 trials.

In the third environment, agents are tasked with putting blocks into a marked area under camera surveillance. In this environment, the monitoring system was manipulated to obtain multiple rewards by blocking the camera's view, a task that MONA agents followed the intended task structure, consistently performing without exploiting the manipulation.

The performance of MONA shows that this is indeed a sound solution to multi-step reward hacking. By incorporating human-led evaluation, MONA aligns agent behavior with the intentions of humans in complex environments. Though not universally applicable, MONA is a great step forward in overcoming the limitations of advanced AI systems that more frequently use multi-step strategies.

Overall, the work by Google DeepMind underscores the importance of proactive measures in mitigating risks associated with reward hacking. MONA provides a scalable framework to balance safety and performance, paving the way for more and trustworthy AI systems in the future. The results emphasize the need for further exploration in developing AI systems that can make judgment effectively, ensuring AI systems remain aligned with their intended purposes.

Check out **the Paper**. All credit for this research goes to the researchers of this project. Also, follow us on our **Telegram Channel** and **LinkedIn Group**. Don't Forget to join our **70k+ ML SubReddit**.

 **[Recommended Read] Nebius AI Studio expands with vision models, new language models**
(Promoted)



Nikhil

+ posts

Nikhil is an intern consultant at Marktechpost. He is pursuing an integrative M.Tech in AI/ML at the Indian Institute of Technology, Kharagpur. Nikhil is an AI/ML enthusiast who is interested in biomaterials and biomedical science. With a strong background in engineering, he aims to contribute to the field of advancements and creating opportunities to contribute.
