

Comparative Robustness of Reward Shaping Techniques Against Adversarial Reward Hacking

Assignee Research

June 7, 2026

Abstract

This report synthesises findings from 10 peer-reviewed papers addressing the following research question: What is the comparative robustness of different reward shaping techniques (e.g., potential-based, state-based) against adversarial reward hacking attacks, as measured by changes in HH-RLHF. 14 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.7/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Reward Shaping to Mitigate Reward Hacking in RLHF. Research question: What is the comparative robustness of different reward shaping techniques (e.g., potential-based, state-based) against adversarial reward hacking attacks, as measured by changes in HH-RLHF helpfulness scores under perturbed inputs?.

2 Methodology

Systematic literature search across multiple databases yielded 10 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.7/10.

3 Results

10 papers retrieved. 14 claims extracted; 0 independently verified. Quality review score: 3.7/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

Claim	Verified	Confidence
The winrate measures the policy model’s winning rate against the SFT model, as evaluated by DeepSeek-V3 (DeepSeek-AI, 2023) and MT-Bench (Zheng et al., 2023a), six metrics are utilized, with al	×	0.03
The SFT model is trained for two epochs on chosen responses with a learning rate of $5e-6$.	×	0.02
The reward model, consisting of a linear head appended to the base model, is trained for one epoch with a learning rate	×	0.06
The policy model, initialized as the SFT model, is trained for one epoch with a learning rate of $3e-7$.	×	0.03
The critic model, initialized as the reward model, is trained for one epoch with a learning rate of $5e-6$.	×	0.05
A linear learning rate scheduler is employed for all training procedures, gradually increasing the learning rate from 0	×	0.03
To generate the reward and winrate curves, the policy model is evaluated on the test set at intervals of 0.1 epochs, yie	×	0.02
Increasing the KL penalty coefficient from 0.01 to 0.1 leads to a rise in the winrate curve and a corresponding decline	×	0.02
A similar effect is observed when reducing the reward ceiling (i.e., the maximum reward threshold).	×	0.02
PAR’s functional form closely resembles the Bradley-Terry model of the proxy reward as an Elo score (Elo, 1978).	×	0.05
The sigmoid transformation effectively suppresses both sources of variance (see Figure 2 and Section 3.3).	×	0.02
PAR demonstrates strong robustness by providing a wider and more forgiving window for early stopping.	×	0.09
We conduct experiments on the base model Gemma2-2B (Google, 2024) using two widely used benchmarks.	×	0.07

References

- <http://arxiv.org/abs/2502.18770v5>
- <http://arxiv.org/abs/1902.06705v2>
- <http://arxiv.org/abs/1801.04693v1>