

Reward Shaping vs. Alternative Alignment Techniques in RLHF for GSM8K Robustness

Assignee Research

June 7, 2026

Abstract

This report synthesises findings from 12 peer-reviewed papers addressing the following research question: How does reward shaping in RLHF pipelines compare to other alignment techniques (e.g., iterative human feedback, adversarial training) in terms of reducing reward hacking incidents on GSM8K while. 12 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.1/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Mitigating Reward Hacking in RLHF via Bayesian Non-negative Reward Modeling. Research question: How does reward shaping in RLHF pipelines compare to other alignment techniques (e.g., iterative human feedback, adversarial training) in terms of reducing reward hacking incidents on GSM8K while maintaining final solution accuracy?.

2 Methodology

Systematic literature search across multiple databases yielded 12 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.1/10.

3 Results

12 papers retrieved. 12 claims extracted; 0 independently verified. Quality review score: 3.1/10.

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
BNRM trained on 1K examples matches the performance of the BT baseline trained on significantly larger datasets.	×	0.03
Under a 40% label noise rate, BNRM improves the BT baseline performance by up to 16.7%.	×	0.04
BNRM trained with 40% label noise rivals the performance of BT trained with only 10%-20% noise.	×	0.03
The RLHF evaluation fine-tuned Llama-3.1-8B-Instruct and OpenRLHF-Llama3-8B-SFT using Proximal Policy Optimization (PPO)	×	0.04
On the GSM8K 4-shot benchmark, the 'Ours' method achieved a score of 77.10 with the OpenRLHF-Llama3-8B-SFT backbone, an	×	0.02
On the Hellaswag benchmark, the 'Ours' method achieved a score of 60.68 with the OpenRLHF-Llama3-8B-SFT backbone, a decr	×	0.02
On the IFEval benchmark, the 'Ours' method achieved a score of 78.20 with the Llama3.1-8B-Instruct backbone, an increase	×	0.03
On the Race 3-shot benchmark, the 'Ours' method achieved a score of 83.31 with the Llama3.1-8B-Instruct backbone, an inc	×	0.02
The BNBT-Reward-Llama-3.1-8B model achieved an average RewardBench score of 93.6.	×	0.03
The BNBT-Reward-Llama-3.1-8B model achieved a Chat-Hard score of 89.7 on Reward-Bench.	×	0.03
The average accuracy of the 'Ours' method on the Llama3.1-8B-Instruct backbone across all listed benchmarks is 74.98.	×	0.02
The average accuracy of the 'Ours' method on the OpenRLHF-Llama3-8B-SFT backbone across all listed benchmarks is 62.25.	×	0.02

References

- <http://arxiv.org/abs/2502.18770v5>
- <http://arxiv.org/abs/2602.10623v2>

- <http://arxiv.org/abs/2308.15969v1>