

Multi-Objective Reward Optimization and Q-Shaping for PowerInfer Throughput Across Programming Languages

Assignee Research

May 30, 2026

Abstract

This report synthesises findings from 15 peer-reviewed papers addressing the following research question: What is the inference efficiency impact of multi-objective reward optimization on PowerInfer's throughput when scaling to diverse programming languages beyond Python. Q-shaping is an extension of Q-value initialization and serves as an alternative to reward shaping for incorporating domain knowledge to accelerate agent training, thereby improving sample efficiency by directly shaping Q-values. This approach is both general and robust across. 9 claims were extracted from source literature; 1 was independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.2/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: From Reward Shaping to Q-Shaping: Achieving Unbiased Learning with LLM-Guided Knowledge. Research question: What is the inference efficiency impact of multi-objective reward optimization on PowerInfer's throughput when scaling to diverse programming languages beyond Python?.

2 Methodology

Systematic literature search across multiple databases yielded 15 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.2/10.

3 Results

15 papers retrieved. 9 claims extracted; 1 independently verified. Quality review score: 4.2/10.

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
Q-shaping improves performance by 16.87% on average compared to the best baseline across 20 tasks.	×	0.14
Q-shaping shows a 55.39% average improvement over TD3 across 20 tasks.	×	0.07
Q-shaping outperforms LLM-based reward shaping methods by 253.80% on average in peak performance.	✓	0.22
LLM-TD3 achieved a 38.68% improvement in the door-close task over the best baseline.	×	0.06
LLM-TD3 achieved a 406.04% improvement in the drawer-open task over the best baseline.	×	0.06
LLM-TD3 achieved a 389.77% improvement in the window-close task over the best baseline.	×	0.06
LLM-TD3 achieved a 180.70% improvement in the sweep-into task over the best baseline.	×	0.06
Most LLMs, including o1-Preview, GPT-4o, DeepSeek-V2.5, and yi-large, provided correct heuristic functions with a 100% s	×	0.05
Gemini achieved only 44% average correctness in providing heuristic functions.	×	0.03

References

- <http://arxiv.org/abs/2410.01458v1>
- <http://arxiv.org/abs/2204.07167v2>

- <http://arxiv.org/abs/2602.20945v3>