

Directional Preference Alignment Robustness to Adversarial Inputs in Code Generation

Assignee Research

May 30, 2026

Abstract

This report synthesises findings from 1 peer-reviewed paper addressing the following research question: How robust is the Directional Preference Alignment framework to adversarial or edge-case inputs in code generation tasks compared to RLHF, as measured by accuracy on a curated subset of HumanEval. The remarkable success of Large Language Models (LLMs) has illuminated a promising pathway toward achieving Artificial General Intelligence for both academic and industrial communities, owing to their unprecedented performance across various applications. As LLMs continue to evolve, 5 claims were extracted from source literature; 5 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 8.0/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: A Comprehensive Survey in LLM(-Agent) Full Stack Safety: Data, Training and Deployment. Research question: How robust is the Directional Preference Alignment framework to adversarial or edge-case inputs in code generation tasks compared to RLHF, as measured by accuracy on a curated subset of HumanEval with perturbed or ambiguous prompts?.

2 Methodology

Systematic literature search across multiple databases yielded 1 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 8.0/10.

3 Results

1 papers retrieved. 5 claims extracted; 5 independently verified. Quality review score: 8.0/10.

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
Existing surveys on LLM safety primarily focus on specific stages of the LLM lifecycle, such as the deployment phase or	✓	0.33
This paper introduces the concept of 'full-stack' safety to systematically consider safety issues throughout the entire	✓	0.35
The paper defines the complete LLM lifecycle as encompassing data preparation, pre-training, post-training, deployment,	✓	0.29
This work represents the first safety survey to encompass the entire lifecycle of LLMs.	✓	0.24
The research is grounded in an exhaustive review of over 800 papers.	✓	0.19

References

- <https://doi.org/10.48550/arxiv.2504.15585>