

Llama-3 and Vicuna Alignment Strategies for Spurious Feature Robustness in Adversarial NLP Benchmarks

Assignee Research

June 8, 2026

Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: How does the alignment strategy in Llama-3 compare to Vicuna’s SFT approach in mitigating spurious feature sensitivity when evaluated on adversarial natural language benchmarks like AdvGLUE or StressTest?. 0 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 2.8/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Assessing Robustness to Spurious Correlations in Post-Training Language Models. Research question: How does the alignment strategy in Llama-3 compare to Vicuna’s SFT approach in mitigating spurious feature sensitivity when evaluated on adversarial natural language benchmarks like AdvGLUE or StressTest?.

2 Methodology

Systematic literature search across multiple databases yielded 13 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 2.8/10.

3 Results

13 papers retrieved. 0 claims extracted; 0 independently verified. Quality review score: 2.8/10.

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.

References

- <http://arxiv.org/abs/2306.11066v2>
- <http://arxiv.org/abs/2505.05704v1>
- <http://arxiv.org/abs/2110.06500v2>