

# Systematic Review of Open-Source Language Model Benchmark Leaderboards

Assignee Research

June 3, 2026

## Abstract

This report synthesises findings from 9 peer-reviewed papers addressing the following research question: Open source language model benchmark leaderboard systematic review. 0 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.0/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Structured Prompts Improve Evaluation of Language Models. Research question: Open source language model benchmark leaderboard systematic review.

## 2 Methodology

Systematic literature search across multiple databases yielded 9 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 5.0/10.

## 3 Results

9 papers retrieved. 0 claims extracted; 0 independently verified. Quality review score: 5.0/10.

## 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.

## References

- <http://arxiv.org/abs/2511.20836v3>
- <http://arxiv.org/abs/2605.30018v2>
- <http://arxiv.org/abs/2507.08538v1>