

Automated Test Suite Robustness for Mistral-Large-2 Code on LiveCodeBench

Assignee Research

May 31, 2026

Abstract

This report synthesises findings from 5 peer-reviewed papers addressing the following research question: What is the robustness of automated test suite evaluations for code generated by Mistral-Large-2 on MBPP when benchmarked against human evaluations using Cohen’s kappa for inter-rater agreement. Large Language Models (LLMs) applied to code-related applications have emerged as a prominent field, attracting significant interest from both academia and industry. However, as new and improved LLMs are developed, existing evaluation benchmarks (e.g., HumanEval, MBPP) are no longer claims were extracted from source literature; 7 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 8.7/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: LiveCodeBench: Holistic and Contamination Free Evaluation of Large Language Models for Code. Research question: What is the robustness of automated test suite evaluations for code generated by Mistral-Large-2 on MBPP when benchmarked against human evaluations using Cohen’s kappa for inter-rater agreement?.

2 Methodology

Systematic literature search across multiple databases yielded 5 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 8.7/10.

3 Results

5 papers retrieved. 7 claims extracted; 7 independently verified. Quality review score: 8.7/10.

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 |
|--|----------|------------|
| LiveCodeBench continuously collects new problems from LeetCode, AtCoder, and CodeForces. | ✓ | 0.22 |
| LiveCodeBench evaluates capabilities including self-repair, code execution, and test output prediction in addition to co | ✓ | 0.22 |
| LiveCodeBench hosts 400 high-quality coding problems published between May 2023 and May 2024. | ✓ | 0.21 |
| The study evaluated 18 base LLMs and 34 instruction-tuned LLMs on LiveCodeBench. | ✓ | 0.24 |
| Existing benchmarks such as HumanEval and MBPP are insufficient for assessing the capabilities of new and improved LLMs. | ✓ | 0.22 |
| The authors will release all prompts and model completions for community analysis. | ✓ | 0.19 |
| The authors will release a general toolkit for adding new scenarios and models. | ✓ | 0.17 |

References

- <https://openalex.org/W7162605369>
- <https://doi.org/10.18653/v1/2025.acl-long.1418>
- <https://doi.org/10.48550/arxiv.2403.07974>