

SOVEREIGN: HumanEval benchmark pass@1 scores comparison GPT-4o Claude Gemini LLaMA DeepSeek 2024

SOVEREIGN Research Kernel
Autonomous draft — Owner review required before publication

May 29, 2026

Abstract

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 sufficient for assessing their capabilities. In this work, we propose LiveCodeBench, a comprehensive and contamination-free evaluation of LLMs for code, which continuously collects new problems over time from contests across three competition platforms, namely LeetCode, AtCoder, and CodeForces. Notably, our benchma

1 Introduction

Analysis of: LiveCodeBench: Holistic and Contamination Free Evaluation of Large Language Models for Code. Research goal: HumanEval benchmark pass@1 scores comparison GPT-4o Claude Gemini LLaMA DeepSeek 2024.

2 Methodology

Multi-query arXiv search (4 parallel queries, Relevance-sorted). TF-IDF cosine semantic verification (bigrams, threshold=0.15). NIM nv-embedqa-e5-v5 (dim=1024) for semantic indexing. Tribunal v2: 3-role parallel review (SKEPTIC/VALIDATOR/SYNTHESIZER) with revision round if score < 6.5.

3 Results

11 papers retrieved. 8 claims extracted, 8 verified. Tribunal: 9.2/10 \$\rightarrow\$ APPROVE (revision_round=0). Policy: AUTO_APPROVE.

4 Uncertainties

NIM free tier latency varies. TF-IDF verification is a weak signal. arXiv Relevance ranking is query-dependent. Tribunal consensus is LLM-based and prompt-sensitive.

5 Extracted Claims

Claim	Verified	Confidence
LiveCodeBench is a comprehensive and contamination-free evaluation benchmark for Large Language Models (LLMs) for code.	✓	0.30
LiveCodeBench continuously collects new problems over time from contests across three competition platforms: LeetCode, A	✓	0.31
LiveCodeBench focuses on a broader range of code-related capabilities, such as self-repair, code execution, and test out	✓	0.36
LiveCodeBench currently hosts four hundred high-quality coding problems that were published between May 2023 and May 202	✓	0.25
18 base LLMs and 34 instruction-tuned LLMs have been evaluated on LiveCodeBench.	✓	0.23
The paper presents empirical findings on contamination, holistic performance comparisons, potential overfitting in exist	✓	0.32
All prompts and model completions from LiveCodeBench will be released for further community analysis.	✓	0.16
A general toolkit for adding new scenarios and models will be released along with LiveCodeBench.	✓	0.18

References

- <https://doi.org/10.48550/arxiv.2406.00515>
- <https://doi.org/10.4230/oasics.icpec.2025.4>
- <https://doi.org/10.48550/arxiv.2403.07974>