

WizardCoder vs. CodeLlama on Multi-Language Code Generation Benchmarks

Assignee Research

June 7, 2026

Abstract

This report synthesises findings from 14 peer-reviewed papers addressing the following research question: How does the performance of WizardCoder on HumanEval compare to other open-source instruction-tuned code LLMs like CodeLlama when evaluated on multi-language code generation tasks, measured by pass@1. 14 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.7/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: WizardCoder: Empowering Code Large Language Models with Evol-Instruct. Research question: How does the performance of WizardCoder on HumanEval compare to other open-source instruction-tuned code LLMs like CodeLlama when evaluated on multi-language code generation tasks, measured by pass@1 accuracy?.

2 Methodology

Systematic literature search across multiple databases yielded 14 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.7/10.

3 Results

14 papers retrieved. 14 claims extracted; 0 independently verified. Quality review score: 3.7/10.

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
HumanEval comprises 164 problems with an average of 9.6 test cases per problem.	×	0.04
HumanEval+ expands the test cases significantly to an average of 774.8 per problem.	×	0.04
MBPP provides 500 test programming problems with three automated test cases each.	×	0.04
WizardCoder-15B achieved a pass@1 score of 57.3% on HumanEval.	×	0.04
WizardCoder-15B achieved a pass@1 score of 51.8% on MBPP.	×	0.03
WizardCoder-34B achieved a pass@1 score of 71.5% on HumanEval.	×	0.04
WizardCoder-34B achieved a pass@1 score of 61.2% on MBPP.	×	0.03
GPT-4 achieved a pass@1 score of 67.0% on HumanEval.	×	0.02
CodeLlama-Python (34B) achieved a pass@1 score of 53.7% on HumanEval.	×	0.05
WizardCoder models demonstrated superior performance across all 8 evaluated programming languages (Java, JavaScript, C++	×	0.10
The DS-1000 benchmark comprises 1k distinct data science workflows spanning 7 libraries.	×	0.04
WizardCoder demonstrates significant superiority over all other models on the DS-1000 benchmark insertion scores.	×	0.06
The pass@1 scores for WizardCoder models were estimated using n=200 samples with temperature=0.2 and top_p=0.95.	×	0.03
MultiPL-E benchmark evaluations used hyperparameters: temperature=0.2, top_p=0.95, max_length=512, and n=50.	×	0.02

References

- <http://arxiv.org/abs/2410.12381v3>
- <http://arxiv.org/abs/2306.08568v2>
- <http://arxiv.org/abs/2403.03788v1>