

Scaling Effects on Instruction Following Accuracy in Out-of-Domain Code Generation

Assignee Research

May 30, 2026

Abstract

This report synthesises findings from 14 peer-reviewed papers addressing the following research question: What is the impact of model scaling on instruction following accuracy when evaluated on out-of-domain code generation tasks. Despite widespread deployment of Large Language Models, systematic evaluation of instruction-following capabilities remains challenging. While comprehensive benchmarks exist, focused assessments that quickly diagnose specific instruction adherence patterns are valuable. 19 claims were extracted from source literature; 2 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.2/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: When Models Can't Follow: Testing Instruction Adherence Across 256 LLMs. Research question: What is the impact of model scaling on instruction following accuracy when evaluated on out-of-domain code generation tasks.

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 4.2/10.

3 Results

14 papers retrieved. 19 claims extracted; 2 independently verified. Quality review score: 4.2/10.

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
The evaluation framework was implemented using the OpenRouter API, which provides unified access to multiple language models	×	0.06
The evaluation system executes each test systematically across all target models, maintaining consistent parameters to ensure reproducibility	×	0.04
Temperature was set to 0.0 to minimize response variability and enable more deterministic evaluation of instruction-following capabilities	×	0.08
Response timeout limits were established at 10 seconds to balance comprehensive model coverage with practical execution	×	0.05
Primary metrics include binary pass/fail determination based on adherence to specified instructions, response time measurement, and token usage	×	0.04
Results are automatically compiled into a structured Excel workbook with multiple analytical views.	×	0.02
The overview sheet provides a high-level pass/fail matrix across all models and tests, enabling rapid identification of model performance trends	×	0.02
Individual test sheets contain detailed response data, allowing deeper investigation of specific failure modes.	×	0.02
A model summary sheet aggregates performance statistics, including success rates, average response times, and total tokens used	×	0.03
Each test prompt includes programmatically verifiable success criteria.	×	0.04
The verification approach accounts for common variations in model outputs while maintaining objective assessment standards	×	0.03
The verification process operates in two stages: primary verification applies strict matching criteria to determine exact matches, while secondary verification uses fuzzy matching for partial success	×	0.02
The comprehensive evaluation of instruction-following capabilities encompassed all 331 models available via OpenRouter at the time of testing	✓	0.19
331 models were verified for basic functionality, of which 256 passed verification and were subsequently evaluated using more complex prompts	×	0.10
The two-stage verification protocol first assessed basic endpoint functionality and subsequently evaluated all verified models against complex instruction-following prompts	×	0.09
The evaluation of instruction-following capabilities was conducted on October 14, 2025.	✓	0.21
The evaluation included 331 models across diverse providers and architectures.	×	0.06
The verification step used the query: 'What is the capital of France?'	×	0.02

References

- <http://arxiv.org/abs/2306.08568v2>
- <http://arxiv.org/abs/2403.09832v1>
- <http://arxiv.org/abs/2510.18892v1>