

# To what extent does instruction complexity in BigCodeBench correlate with the performance degradation of code

Assignee Research

May 29, 2026

## **Abstract**

Large Language Models (LLMs) have garnered remarkable advancements across diverse code-related tasks, known as Code LLMs, particularly in code generation that generates source code with LLM from natural language descriptions. This burgeoning field has captured significant interest from both academic researchers and industry professionals due to its practical significance in software development, e.g., GitHub Copilot. Despite the active exploration of LLMs for a variety of code tasks, either from the perspective of natural language processing (NLP) or software engineering (SE) or both, there is

## **1 Introduction**

This paper examines: A Survey on Large Language Models for Code Generation. Research question: To what extent does instruction complexity in BigCodeBench correlate with the performance degradation of code generation models across different Python data science libraries?.

## **2 Methodology**

Systematic literature search across multiple databases yielded 13 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.4/10.

## **3 Results**

13 papers retrieved. 6 claims extracted; 5 independently verified. Quality review score: 7.4/10.

## 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
Large Language Models (LLMs) have achieved advancements in code generation tasks, generating source code from natural la	✓	0.17
GitHub Copilot is an example of a practical application of LLMs for code generation in software development.	×	0.14
There is a noticeable absence of a comprehensive and up-to-date literature review dedicated specifically to LLMs for cod	✓	0.25
The survey introduces a taxonomy categorizing developments in LLMs for code generation covering data curation, latest ad	✓	0.31
The survey presents an empirical comparison of LLMs using the HumanEval, MBPP, and Big-CodeBench benchmarks.	✓	0.16
The empirical comparison in the survey covers various levels of difficulty and types of programming tasks.	✓	0.19

## References

- <https://doi.org/10.54254/2755-2721/2025.tj22242>
- <https://doi.org/10.48550/arxiv.2406.15877>
- <https://doi.org/10.48550/arxiv.2406.00515>