

# Robustness Degradation in Self-Invoking Code Generation under Adversarial Perturbations

Assignee Research

June 8, 2026

## Abstract

This report synthesises findings from 14 peer-reviewed papers addressing the following research question: How does the robustness of self-invoking code generation on MBPP Pro degrade under adversarial perturbations in the base problem description across different model families. 0 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.0/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: HumanEval Pro and MBPP Pro: Evaluating Large Language Models on Self-invoking Code Generation. Research question: How does the robustness of self-invoking code generation on MBPP Pro degrade under adversarial perturbations in the base problem description across different model families?.

## 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 3.0/10.

## 3 Results

14 papers retrieved. 0 claims extracted; 0 independently verified. Quality review score: 3.0/10.

## 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.

## References

- <http://arxiv.org/abs/2110.09903v2>
- <http://arxiv.org/abs/2008.07651v1>
- <http://arxiv.org/abs/2412.21199v2>