

RLHF Alignment and Adversarial Robustness in Sparse MoE Code Generation Models

Assignee Research

May 30, 2026

Abstract

This report synthesises findings from 16 peer-reviewed papers addressing the following research question: How does RLHF alignment impact the adversarial robustness of sparse MoE models on code generation benchmarks like HumanEval Pro compared to dense architectures. We introduce self-invoking code generation, a new task designed to evaluate the progressive reasoning and problem-solving capabilities of LLMs. In this task, models are presented with a base problem and a related, more complex problem. 15 claims were extracted from source literature; 3 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.5/10. 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 RLHF alignment impact the adversarial robustness of sparse MoE models on code generation benchmarks like HumanEval Pro compared to dense architectures?.

2 Methodology

Systematic literature search across multiple databases yielded 16 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.5/10.

3 Results

16 papers retrieved. 15 claims extracted; 3 independently verified. Quality review score: 4.5/10.

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
o1-mini achieves 96.2% pass@1 on HumanEval.	✓	0.20
o1-mini achieves 76.2% pass@1 on HumanEval Pro.	✓	0.22
Instruction-tuned models are less efficient on self-invoking code generation than on traditional code generation tasks.	✓	0.29
HumanEval and MBPP are fundamental benchmarks focusing on Python function completion tasks with test-driven evaluation.	×	0.10
Deepseek-V2.5 was used to generate self-invoking problems, candidate solutions, and test inputs for the benchmark.	×	0.03
OpenCoder-8B-base achieved a score of 56.1 on the Base Problem and 10.5 on the Self-invoking Problem.	×	0.08
OpenCoder-8B-instruct achieved a score of 75.4 on the Base Problem and 22.8 on the Self-invoking Problem.	×	0.07
DeepseekCoder-6.7B-base achieved a score of 59.6 on the Base Problem and 35.1 on the Self-invoking Problem.	×	0.08
DeepseekCoder-6.7B-instruct achieved a score of 56.1 on the Base Problem and 35.1 on the Self-invoking Problem.	×	0.08
WaveCoder-Ultra-6.7B achieved a score of 61.4 on the Base Problem and 26.3 on the Self-invoking Problem.	×	0.07
Magocoder-S-DS-6.7B achieved a score of 50.9 on the Base Problem and 33.3 on the Self-invoking Problem.	×	0.07
Yi-Coder-9B-Chat achieved a score of 66.7 on the Base Problem and 31.6 on the Self-invoking Problem.	×	0.06
Qwen2.5Coder-7B-base achieved a score of 59.6 on the Base Problem and 38.6 on the Self-invoking Problem.	×	0.07
Qwen2.5Coder-7B-instruct achieved a score of 64.9 on the Base Problem and 35.1 on the Self-invoking Problem.	×	0.07
DeepseekCoder-33B-instruct achieved a score of 80.7 on the Base Problem and 43.9 on the Self-invoking Problem.	×	0.07

References

- <http://arxiv.org/abs/2412.15004v4>
- <http://arxiv.org/abs/2307.02055v1>
- <http://arxiv.org/abs/2412.21199v2>