

Multimodal vs. Text-Only Code Models Throughput on HumanEval-V Benchmark

Assignee Research

May 31, 2026

Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: What is the throughput comparison of multimodal code models with visual encoders versus text-only models when evaluated on the HumanEval-V benchmark. Understanding and reasoning over diagrams is a fundamental aspect of human intelligence. While Large Multimodal Models (LMMs) have demonstrated impressive capabilities across various tasks, existing benchmarks lack comprehensive evaluation of their diagram interpretation and. 12 claims were extracted from source literature; 2 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.9/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: HumanEval-V: Benchmarking High-Level Visual Reasoning with Complex Diagrams in Coding Tasks. Research question: What is the throughput comparison of multimodal code models with visual encoders versus text-only models when evaluated on the HumanEval-V benchmark?.

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 4.9/10.

3 Results

13 papers retrieved. 12 claims extracted; 2 independently verified. Quality review score: 4.9/10.

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-V consists of 253 human-annotated coding tasks.	✓	0.16
Each task in HumanEval-V features a diagram encoding the problem context, a function signature defining the task’s input	×	0.08
The top-performing model, Claude 3.5 Sonnet, achieves 36.8% pass@1 on HumanEval-V.	×	0.10
The best open-weight model, Pixtral 124B, reaches 21.3% pass@1 on HumanEval-V.	×	0.03
Claude 3.5 Sonnet achieves a 74.3% pass rate with 100 samples on HumanEval-V.	×	0.05
Claude 3.5 Sonnet can reach 55.3% pass@1 with four self-refining iterations based on test case execution feedback on Hum	×	0.04
HumanEval-V offers a more diverse and complex set of diagrams spanning six task types.	×	0.11
HumanEval-V demands versatile capabilities for diagram understanding and reasoning.	×	0.11
The visual context must be essential for solving the task in HumanEval-V, with all relevant information contained in a s	×	0.05
Tasks in HumanEval-V should be designed around the visual context with minimal textual description.	×	0.07
The two-stage evaluation pipeline in HumanEval-V supports LMMs with limited coding abilities by first prompting them to	×	0.07
Extensive experiments with 22 LMMs were conducted on HumanEval-V.	✓	0.15

References

- <http://arxiv.org/abs/2410.12381v3>
- <http://arxiv.org/abs/2412.09616v2>
- <http://arxiv.org/abs/2407.04973v1>