

# Multimodal vs. Text-Only Models on HumanEval-V Diagram-Based Reasoning Accuracy

Assignee Research

May 31, 2026

## Abstract

This report synthesises findings from 4 peer-reviewed papers addressing the following research question: How do multimodal models perform on HumanEval-V compared to text-only models when evaluated with accuracy metrics on diagram-based reasoning tasks. 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. 16 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 2.7/10. 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: How do multimodal models perform on HumanEval-V compared to text-only models when evaluated with accuracy metrics on diagram-based reasoning tasks?.

## 2 Methodology

Systematic literature search across multiple databases yielded 4 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 2.7/10.

### **3 Results**

4 papers retrieved. 16 claims extracted; 0 independently verified. Quality review score: 2.7/10.

### **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.	×	0.15
Each task in HumanEval-V features a diagram encoding the problem context, a function signature defining the task’s input	×	0.07
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.	×	0.02
Claude 3.5 Sonnet achieves a 74.3% pass rate with 100 samples on HumanEval-V.	×	0.04
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.	×	0.10
HumanEval-V demands versatile capabilities for diagram understanding and reasoning.	×	0.11
HumanEval-V uses code generation tasks for evaluation instead of multiple-choice or short-answer questions.	×	0.08
HumanEval-V addresses the lack of benchmarks for evaluating visual reasoning in programming contexts.	×	0.12
HumanEval-V’s tasks are designed around the visual context with minimal textual description.	×	0.04
HumanEval-V’s visual context must be essential for solving the task, with all relevant information contained in a single	×	0.04
HumanEval-V’s test cases rigorously verify whether the model captures all critical visual information.	×	0.05
HumanEval-V utilizes a two-stage evaluation pipeline that supports LMMs with limited coding abilities.	×	0.08
HumanEval-V’s evaluation pipeline prioritizes visual understanding over coding proficiency.	×	0.06
HumanEval-V was evaluated with 22 LMMs.	×	0.10

## References

- <http://arxiv.org/abs/2410.12381v3>
- <http://arxiv.org/abs/2406.12934v1>
- <http://arxiv.org/abs/2411.13760v1>