

State-of-the-art LMMs vs Specialized Models in Diagram Understanding on HumanEval-V

Assignee Research

June 8, 2026

Abstract

This report synthesises findings from 14 peer-reviewed papers addressing the following research question: What is the performance gap between state-of-the-art LMMs and specialized diagram understanding models on HumanEval-V, and how does this gap vary by task type. 14 claims were extracted from source literature; 1 was 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-V: Benchmarking High-Level Visual Reasoning with Complex Diagrams in Coding Tasks. Research question: What is the performance gap between state-of-the-art LMMs and specialized diagram understanding models on HumanEval-V, and how does this gap vary by task type?.

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 4.5/10.

3 Results

14 papers retrieved. 14 claims extracted; 1 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
HumanEval-V consists of 253 human-annotated coding tasks.	✓	0.16
Each task in HumanEval-V features a diagram, a function signature, and test cases.	×	0.12
HumanEval-V diagrams span six task types.	×	0.12
HumanEval-V uses code generation tasks for evaluation instead of multiple-choice or short-answer questions.	×	0.09
Claude 3.5 Sonnet achieves a 36.8% pass@1 score on HumanEval-V.	×	0.11
Pixtral 124B achieves a 21.3% pass@1 score on HumanEval-V.	×	0.03
Claude 3.5 Sonnet achieves a 74.3% pass rate with 100 samples.	×	0.04
Claude 3.5 Sonnet reaches a 55.3% pass@1 rate with four self-refining iterations based on test case execution feedback.	×	0.04
The evaluation pipeline involves prompting models to generate a structured diagram description before using a coder mode	×	0.05
Experiments were conducted with 22 Large Multimodal Models (LMMs).	×	0.13
GPT-4o achieved a 27.7% pass@1 score in the baseline setting according to Table (p5).	×	0.03
Gemini 1.5 Pro achieved a 22.9% pass@1 score in the baseline setting according to Table (p5).	×	0.02
Pixtral 124B achieved a 16.6% pass@1 score in the baseline setting according to Table (p5).	×	0.01
The Intermediate Textual Representation method requires the model to produce a structured textual problem specification	×	0.06

References

- <http://arxiv.org/abs/2311.02782v3>
- <http://arxiv.org/abs/2508.19294v2>
- <http://arxiv.org/abs/2410.12381v3>