

# Vision-Language Models vs. Text-Only LLMs on HumanEval-V with Chain-of-Thought Prompting

Assignee Research

June 6, 2026

## Abstract

This report synthesises findings from 11 peer-reviewed papers addressing the following research question: How do vision-language models compare to text-only LLMs in accuracy on HumanEval-V when evaluated with chain-of-thought prompting. 9 claims were extracted from source literature; 9 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 9.0/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Large language models encode clinical knowledge. Research question: How do vision-language models compare to text-only LLMs in accuracy on HumanEval-V when evaluated with chain-of-thought prompting?.

## 2 Methodology

Systematic literature search across multiple databases yielded 11 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 9.0/10.

## 3 Results

11 papers retrieved. 9 claims extracted; 9 independently verified. Quality review score: 9.0/10.

## 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.

## 5 Extracted Claims

Claim	Verified	Confidence
Large language models (LLMs) have demonstrated impressive capabilities, but the bar for clinical applications is high.	✓	0.25
Attempts to assess the clinical knowledge of models typically rely on automated evaluations based on limited benchmarks.	✓	0.29
MultiMedQA is a benchmark combining six existing medical question answering datasets spanning professional medicine, res	✓	0.38
A human evaluation framework for model answers along multiple axes including factuality, comprehension, reasoning, possi	✓	0.30
Flan-PaLM achieves state-of-the-art accuracy on every MultiMedQA multiple-choice dataset (MedQA, MedMCQA, PubMedQA and M	✓	0.36
Flan-PaLM achieves 67.6% accuracy on MedQA (US Medical Licensing Exam-style questions), surpassing the prior state of th	✓	0.34
Human evaluation reveals key gaps in the performance of Flan-PaLM.	✓	0.20
Instruction prompt tuning is a parameter-efficient approach for aligning LLMs to new domains using a few exemplars.	✓	0.29
Med-PaLM performs encouragingly, but remains inferior to clinicians.	✓	0.21

## References

- <https://doi.org/10.1038/s41586-023-06291-2>
- <https://doi.org/10.1007/s11704-026-60308-3>
- <https://doi.org/10.18653/v1/2023.findings-acl.29>