

Dense Annotation Pipelines Enhance Multimodal Model Robustness in Low-Light and Occluded Environments

Assignee Research

June 7, 2026

Abstract

This report synthesises findings from 16 peer-reviewed papers addressing the following research question: What is the effect of dense annotation pipelines (e.g., EPIC-KITCHENS-100's 54% increase) on the robustness of multimodal models in low-light or occluded environments, as measured by accuracy on. 11 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.3/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Perception, Reason, Think, and Plan: A Survey on Large Multimodal Reasoning Models. Research question: What is the effect of dense annotation pipelines (e.g., EPIC-KITCHENS-100's 54% increase) on the robustness of multimodal models in low-light or occluded environments, as measured by accuracy on Ego4D benchmarks?.

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 3.3/10.

3 Results

16 papers retrieved. 11 claims extracted; 0 independently verified. Quality review score: 3.3/10.

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
LLaVA-CoT (2024a) introduces LLaVA-CoT-100k and scalable beam search.	×	0.01
LlamaV-o1 (2025) introduces VCR-Bench and outperforms previous models.	×	0.01
Mulberry (2024a) introduces Mulberry-260k and CoMCTS for collective learning.	×	0.01
RedStar-Geo (2025a) is competitive with minimal Long-CoT data.	×	0.03
RBF++ (2025d) proposes SR-MCTS for structured search and PPRM for evaluating reasoning boundaries.	×	0.03
Multimodal-O1 models extend System-1 reasoning by deepening CoT workflows through multi-stage generation structures, lon	×	0.05
Multimodal-O1 models achieve more coherent, interpretable, and scalable multimodal reasoning.	×	0.11
Early multimodal models primarily focused on the representation, alignment, and fusion of information.	×	0.09
Reasoning in early multimodal models was often implicit, typically requiring separate, task-specific reasoning modules.	×	0.10
Multimodal large language models, particularly those adopting a vision encoder-language model structure, have achieved a	×	0.09
Multimodal large language models have demonstrated improved multi-task reasoning performance.	×	0.08

References

- <http://arxiv.org/abs/2205.06218v1>
- <http://arxiv.org/abs/2303.06705v3>
- <http://arxiv.org/abs/2505.04921v2>