

# Scaling Unimodal Instruction Data for Multimodal CLAM Performance in RoboSuite

Assignee Research

June 9, 2026

## Abstract

This report synthesises findings from 4 peer-reviewed papers addressing the following research question: How does scaling the amount of unimodal instruction data during training impact the multimodal task accuracy of CLAM in RoboSuite, when benchmarked against vision-only and audio-only baselines using. 13 claims were extracted from source literature; 11 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 7.5/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Multimodal Large Language Models for Diagnostic Feedback Analytics in STEM Learning Platforms. Research question: How does scaling the amount of unimodal instruction data during training impact the multimodal task accuracy of CLAM in RoboSuite, when benchmarked against vision-only and audio-only baselines using the RoboSuite success rate metric?.

## 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 7.5/10.

## 3 Results

4 papers retrieved. 13 claims extracted; 11 independently verified. Quality review score: 7.5/10.

## 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
Most existing automated feedback systems rely on unimodal inputs or rule-based logic.	✓	0.21
Existing unimodal or rule-based feedback systems provide surface-level feedback that fails to capture underlying learner	✓	0.26
The proposed framework integrates heterogeneous learner data including text responses, symbolic mathematics, diagrams, c	✓	0.29
The framework uses modality-specific encoders and attention-based fusion strategies to integrate learner data.	✓	0.18
Diagnostic reasoning in the framework is performed using a multimodal large language model constrained by curricular obj	✓	0.30
The framework includes explainability mechanisms such as rationale tracing and attention visualization.	✓	0.16
Empirical evaluation was conducted across mathematics, physics, and computer science tasks.	✓	0.18
The proposed framework demonstrates significant improvements over baseline systems in diagnostic accuracy.	✓	0.17
The proposed framework demonstrates significant improvements over baseline systems in learning gains.	✓	0.16
The proposed framework demonstrates significant improvements over baseline systems in error correction rates.	✓	0.15
The proposed framework demonstrates significant improvements over baseline systems in learner engagement.	×	0.14
The proposed framework demonstrates significant improvements over baseline systems in trust.	×	0.11
Multimodal LLM-driven diagnostic feedback can operationalize formative assessment principles at scale.	✓	0.30

## References

- <https://arxiv.org/abs/2506.01111>
- <https://www.semanticscholar.org/paper/f9d0c3031bdf285ece911f49d6c27225f91f6ace>
- <https://www.semanticscholar.org/paper/78fb640068f3b6e881ccbc7f005460f8ebf7545b>