

Multimodal Language Models in Visual Mathematical and Scientific Reasoning Benchmarks

Assignee Research

June 7, 2026

Abstract

This report synthesises findings from 16 peer-reviewed papers addressing the following research question: How do multimodal language models perform on visual mathematical and scientific reasoning v18. 15 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: Hierarchical Pre-Training of Vision Encoders with Large Language Models. Research question: How do multimodal language models perform on visual mathematical and scientific reasoning v18.

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 4.5/10.

3 Results

16 papers retrieved. 15 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
HIVE is evaluated on CIFAR-10, CIFAR-100, ImageNet-1K, Tiny-ImageNet, Food-101, Stanford Cars, Oxford-IIIT Pets, and Cal	×	0.02
HIVE is evaluated on MME, GQA, OK-VQA, and ScienceQA for vision-language model (VLM) evaluation.	✓	0.17
HIVE is compared against two baseline configurations: Base and SA.	×	0.02
Base configuration uses the original foundation models CLIP and SigLIP without additional LLM-supported pre-training.	×	0.06
SA configuration is a self-attention-based vision encoder trained using a three-stage pre-training method.	×	0.15
Both SA and HIVE follow identical pre-training strategies for the vision encoder.	×	0.12
For VLM evaluation, both SA and HIVE are further fine-tuned following the procedure used in LLaVA.	×	0.05
MobileLLM-350M is used as the language model backbone for both self-attention and hierarchical cross-attention configura	×	0.11
Optimization is performed using decoupled AdamW with a peak learning rate of 1×10^{-3} and a cosine decay schedule.	×	0.02
Gradient clipping and a linear warmup phase are applied to maintain stability.	×	0.02
For VLM fine-tuning, LLaVA is adopted using the Llama-3.2-1B-Instruct model.	×	0.02
For classification tasks, a classifier head is appended to the vision encoder and fine-tuned while keeping the vision en	×	0.08
All experiments are conducted on a single RTX 3090 GPU with a maximum batch size of 256.	×	0.02
HIVE consistently outperforms self-attention baselines across classification and vision-language benchmarks.	×	0.10
HIVE significantly reduces computational overhead.	×	0.02

References

- <http://arxiv.org/abs/2407.04973v1>
- <http://arxiv.org/abs/2604.00086v1>
- <http://arxiv.org/abs/2508.15802v1>