

Graph-Enhanced Multimodal Models vs. Flat Fusion in MM-Vet Throughput-Accuracy Trade-offs

Assignee Research

June 1, 2026

Abstract

This report synthesises findings from 4 peer-reviewed papers addressing the following research question: How do graph-enhanced multimodal models scale in terms of throughput versus accuracy compared to flat fusion models under non-adversarial conditions on MM-Vet. Robot vision has greatly benefited from advancements in multimodal fusion techniques and vision-language models (VLMs). We adopt a task-oriented perspective to systematically review the applications and advancements of multimodal fusion methods and VLMs in the field of robot. 16 claims were extracted from source literature; 2 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.4/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Multimodal Fusion and Vision-Language Models: A Survey for Robot Vision. Research question: How do graph-enhanced multimodal models scale in terms of throughput versus accuracy compared to flat fusion models under non-adversarial conditions on MM-Vet?.

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 4.4/10.

3 Results

4 papers retrieved. 16 claims extracted; 2 independently verified. Quality review score: 4.4/10.

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
Early fusion directly fuses data from different modalities before feature extraction.	×	0.02
Mid-term fusion combines modal features through specific mechanisms such as feature concatenation or weighting after ext	×	0.02
Late stage fusion is achieved by integrating the decision results of each modality after independent decision-making is	×	0.05
Transformer structures have been proposed to improve the applicability of different modal data and capture local feature	×	0.02
Adversarial representation learning is used to create modality invariant embedding spaces, reduce modal gaps, and improv	×	0.06
Post fusion combines the results of decision level independent processing of modalities.	×	0.03
Common techniques in post fusion include weighted averaging, voting mechanisms, and logical rules.	×	0.03
Post fusion offers advantages such as strong modal independence, ease of individual optimization, and scalability of mul	×	0.02
Roitberg et al. compared and analyzed seven decision-level fusion strategies for driver behavior understanding.	×	0.03
Traditional multimodal fusion methods struggle with complex data compared to deep neural networks.	×	0.13
Deep neural networks have made fusion stages less distinct by deeply integrating feature extraction, modality interactio	×	0.03
There has been a shift from explicit to implicit fusion where network design inherently captures modality relationships.	×	0.01
Multimodal fusion approaches in semantic scene understanding are categorized into encoder-decoder frameworks, attention-	✓	0.20
The encoder-decoder method represents scene semantics through encoding, interaction, and decoding.	×	0.04
Various sensory inputs such as RGB, Depth, LiDAR, GPS, and IMU are processed through multimodal fusion strategies to enh 4	×	0.11
Fused features support core robotic vision tasks including 3D semantic scene understanding, SLAM, 3D object detection, n	✓	0.18

References

- <http://arxiv.org/abs/2504.02477v3>
- <http://arxiv.org/abs/2408.00765v2>
- <http://arxiv.org/abs/2308.02490v4>