

How does the semantic-guided diffusion fusion in multimodal VSLAM systems compare to early fusion baselines in

Assignee Research

June 10, 2026

Abstract

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 vision. For semantic scene understanding tasks, we categorize fusion approaches into encoder-decoder frameworks, attention-based architectures, and graph neural networks. Meanwhile, we also analyze the architectural characteristics and practical implementations of these fusion strategies in key tasks such as simultaneous l

1 Introduction

This paper examines: Multimodal Fusion and Vision-Language Models: A Survey for Robot Vision. Research question: How does the semantic-guided diffusion fusion in multimodal VSLAM systems compare to early fusion baselines in terms of feature matching accuracy under high-parallax conditions, measured using the HPatches benchmark?.

2 Methodology

Systematic literature search across multiple databases yielded 7 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

7 papers retrieved. 7 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
The transformer structure has been proposed to improve the applicability of different modal data and capture local featu	×	0.06
Adversarial representation learning has been used to create modality invariant embedding spaces, reduce modal gaps, and	×	0.05
Post fusion is a key method in multimodal analysis, which combines the results of decision level independent processing	×	0.05
Common techniques in post fusion include weighted averaging, voting mechanisms, and logical rules.	×	0.04
Roitberg et al. compared and analyzed seven decision-level fusion strategies for driver behavior understanding.	×	0.03
The encoder-decoder method efficiently represents scene semantics through encoding, interaction, and decoding.	×	0.04
Attention-based fusion has been used in multimodal fusion approaches for semantic scene understanding.	✓	0.18

References

- <http://arxiv.org/abs/2302.04024v1>
- <http://arxiv.org/abs/2504.02477v3>
- <http://arxiv.org/abs/2508.05264v6>