

MobileVLM and State-of-the-Art VLMs on MM1K Under Low-Resource Robotic Manipulation

Assignee Research

June 5, 2026

Abstract

This report synthesises findings from 4 peer-reviewed papers addressing the following research question: What is the performance gap between MobileVLM and state-of-the-art VLMs on the MM1K benchmark when evaluated under low-resource settings (e.g., 5-shot learning) for robotic manipulation tasks. 8 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: Large VLM-based Vision-Language-Action Models for Robotic Manipulation: A Survey. Research question: What is the performance gap between MobileVLM and state-of-the-art VLMs on the MM1K benchmark when evaluated under low-resource settings (e.g., 5-shot learning) for robotic manipulation tasks?.

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 3.3/10.

3 Results

4 papers retrieved. 8 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
Single-system VLA models aim to transfer the semantic knowledge of large VLMs to robotic manipulation tasks through a un	×	0.14
Single-system VLAs offer architectural simplicity, streamlined development, and avoidance of complex inter-module commun	×	0.03
Dual-system VLAs leverage a division of labor to combine reactive speed with deliberate accuracy.	×	0.01
Autoregressive Decoding is a classic paradigm in single-system VLA models.	×	0.10
Inference Efficiency Optimization includes Architectural Optimization, Parameter Optimization, and Inference Acceleratio	×	0.03
Model Performance Enhancement includes Enhancing Perception Modalities and Enhancing Reasoning Capabilities.	×	0.03
Cascade-based Methods include Separate Action Expert and Unified Action Expert.	×	0.06
Parallel-based Methods include Shared-attention Architecture and Cross-attention Architecture.	×	0.03

References

- <http://arxiv.org/abs/2508.13073v2>
- <http://arxiv.org/abs/1911.01557v2>
- <http://arxiv.org/abs/2307.03659v1>