

Gemini 1.5 Flash and LLaVA-NeXT Inference Latency on Video-MME Under 24GB VRAM Constraints

Assignee Research

May 30, 2026

Abstract

This report synthesises findings from 4 peer-reviewed papers addressing the following research question: How does the inference latency of Gemini 1.5 Flash compare to LLaVA-NeXT on the Video-MME benchmark when constrained to 24GB VRAM. In this work, we present a novel method to tackle the token generation challenge in Vision Language Models (VLMs) for video and image understanding, called LLaMA-VID. Current VLMs, while proficient in tasks like image captioning and visual question answering, face computational. 0 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.8/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: LLaMA-VID: An Image is Worth 2 Tokens in Large Language Models. Research question: How does the inference latency of Gemini 1.5 Flash compare to LLaVA-NeXT on the Video-MME benchmark when constrained to 24GB VRAM?.

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 3.8/10.

3 Results

4 papers retrieved. 0 claims extracted; 0 independently verified. Quality review score: 3.8/10.

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.

References

- <https://doi.org/10.18653/v1/2025.findings-emnlp.1127>
- <https://doi.org/10.1007/s00521-024-09426-2>
- <https://doi.org/10.48550/arxiv.2311.17043>