

SOVEREIGN: How does the predictive expert caching latency and token scheduling overhead affect end-to-end tokens-per-second

SOVEREIGN Research Kernel

Autonomous draft — Owner review required before publication

May 29, 2026

Abstract

We present DeepSeek-VL2, an advanced series of large Mixture-of-Experts (MoE) Vision-Language Models that significantly improves upon its predecessor, DeepSeek-VL, through two key major upgrades. For the vision component, we incorporate a dynamic tiling vision encoding strategy designed for processing high-resolution images with different aspect ratios. For the language component, we leverage DeepSeekMoE models with the Multi-head Latent Attention mechanism, which compresses Key-Value cache into latent vectors, to enable efficient inference and high throughput. Trained on an improved vi

1 Introduction

Analysis of: DeepSeek-VL2: Mixture-of-Experts Vision-Language Models for Advanced Multimodal Understanding. Research goal: How does the predictive expert caching latency and token scheduling overhead affect end-to-end tokens-per-second throughput on multimodal reasoning benchmarks (MMMUS, MathVista) for MoE-LLaVA compared to dense model baselines at 7B and 13B parameter scales?.

2 Methodology

Multi-query arXiv search (4 parallel queries, Relevance-sorted). TF-IDF cosine semantic verification (bigrams, threshold=0.15). NIM nv-embedqa-e5-v5 (dim=1024) for semantic indexing. Tribunal v2: 3-role parallel review (SKEPTIC/VALIDATOR/SYNTHESIZER) with revision round if score < 6.5.

3 Results

6 papers retrieved. 6 claims extracted, 6 verified. Tribunal: 7.8/10 \rightarrow APPROVE (revision_round=0). Policy: AUTO_APPROVE.

4 Uncertainties

NIM free tier latency varies. TF-IDF verification is a weak signal. arXiv Relevance ranking is query-dependent. Tribunal consensus is LLM-based and prompt-sensitive.

5 Extracted Claims

Claim	Verified	Confidence
DeepSeek-VL2 incorporates a dynamic tiling vision encoding strategy designed for processing high-resolution images with	✓	0.27
DeepSeek-VL2 leverages DeepSeekMoE models with the Multi-head Latent Attention mechanism, which compresses Key-Value cac	✓	0.27
DeepSeek-VL2 demonstrates superior capabilities across various tasks including visual question answering, optical charac	✓	0.30
DeepSeek-VL2 model series includes three variants: DeepSeek-VL2-Tiny, DeepSeek-VL2-Small, and DeepSeek-VL2, with 1.0B, 2	✓	0.27
DeepSeek-VL2 achieves competitive or state-of-the-art performance with similar or fewer activated parameters compared to	✓	0.30
Codes and pre-trained models are publicly accessible at https://github.com/deepseek-ai/DeepSeek-VL2 .	✓	0.26

References

- <https://doi.org/10.48550/arxiv.2403.14608>
- <https://doi.org/10.55056/jec.1000>
- <https://doi.org/10.48550/arxiv.2412.10302>