

# SOVEREIGN: How does the robustness of SMoES-based MoE-VLMs with soft modality-guided routing compare to dense models of e

SOVEREIGN Research Kernel  
Autonomous draft — Owner review required before publication

May 28, 2026

## Abstract

Mixture-of-Experts (MoE) has become a prevalent backbone for large vision-language models (VLMs), yet how modality-specific signals should guide expert routing remains under-explored. Existing routing strategies are either hand-crafted or modality-agnostic, relying on idealized priors that ignore the layer-dependent modality fusion patterns in MoE-VLMs and provide little guidance for expert specialization. We propose Soft Modality-guided Expert Specialization (SMoES), which consists of dynamic soft modality scores that capture layer-dependent fusion patterns, an expert binning mechanism aligne

## 1 Introduction

Analysis of: SMOES: Soft Modality-Guided Expert Specialization in MoE-VLMs. Research goal: How does the robustness of SMOES-based MoE-VLMs with soft modality-guided routing compare to dense models of equivalent size on adversarial multimodal inputs from the MMMU benchmark at 7B and 13B 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**

11 papers retrieved. 10 claims extracted, 0 verified. Tribunal: 3.3/10 → REJECT (revision\_round=0). Policy: ESCALATE\_TO\_OWNER.

### **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
SMoES achieves a 10.3% reduction in TTFT compared to the baseline on MMMU task with batch size 1.	×	0.02
SMoES achieves a 10.5% reduction in TPOT compared to the baseline on MMMU task with batch size 1.	×	0.02
SMoES achieves a 22.0% reduction in TTFT compared to the baseline on MMMU task with batch size 16.	×	0.03
SMoES achieves a 9.2% reduction in TTFT compared to the baseline on SQA-IMG task with batch size 1.	×	0.02
SMoES achieves a 9.7% reduction in TPOT compared to the baseline on SQA-IMG task with batch size 1.	×	0.01
SMoES achieves a 16.6% reduction in TTFT compared to the baseline on SQA-IMG task with batch size 8.	×	0.02
SMoES achieves a 13.0% reduction in TTFT compared to the baseline on MMMU task with batch size 4.	×	0.02
SMoES achieves a 15.7% reduction in TTFT compared to the baseline on MMMU task with batch size 8.	×	0.02
SMoES achieves a 23.9% reduction in TTFT compared to the baseline on MMMU task with batch size 32.	×	0.04
SMoES achieves a 22.6% reduction in TTFT compared to the baseline on MMMU task with batch size 32.	×	0.04

## References

- <http://arxiv.org/abs/2604.23996v1>
- <http://arxiv.org/abs/2103.15670v3>
- <http://arxiv.org/abs/2603.11114v1>