

# SOVEREIGN: What is the accuracy difference on Winoground between SMOES and modality-agnostic MoE-VLMs when using top-1 ve

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: What is the accuracy difference on Winoground between SMOES and modality-agnostic MoE-VLMs when using top-1 versus top-4 routing under a fixed 64-expert configuration?.

## 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. 7 claims extracted, 0 verified. Tribunal: 3.7/10 → 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 latency on SQA-IMG task at batch size 1	×	0.02
SMoES achieves a 9.2% reduction in TTFT and 9.7% reduction in TPOT on SQA-IMG task	×	0.02
SMoES achieves a 13.0% reduction in TTFT on SQA-IMG task at batch size 16	×	0.02
SMoES achieves a 22.0% reduction in TTFT and 9.0% reduction in TPOT on MMMU task at batch size 16	×	0.03
SMoES achieves a 16.6% reduction in TTFT and 11.3% reduction in TPOT on SQA-IMG task at batch size 16	×	0.03
SMoES achieves 15.7% reduction in TTFT and 9.3% reduction in TPOT on MMMU task at batch size 4	×	0.02
SMoES achieves 23.9% reduction in TTFT and 8.6% reduction in TPOT on MMMU task at batch size 16	×	0.03

## References

- <http://arxiv.org/abs/2604.23996v1>
- <http://arxiv.org/abs/2507.17467v1>
- <http://arxiv.org/abs/2603.11114v1>