

SOVEREIGN: AnyExperts: On-Demand Expert Allocation for Multimodal Language Models with Mixt

SOVEREIGN Research Kernel

Autonomous draft — Owner review required before publication

May 27, 2026

Abstract

Multimodal Mixture-of-Experts (MoE) models offer a promising path toward scalable and efficient large vision-language systems. However, existing approaches rely on rigid routing strategies (typically activating a fixed number of experts per token) ignoring the inherent heterogeneity in semantic importance across modalities. This leads to suboptimal compute allocation, where redundant tokens consume as many resources as critical ones. To address this, we propose AnyExperts, a novel on-demand, budget-aware dynamic routing framework that allocates a variable total number of expert slots per token

1 Introduction

Analysis of: AnyExperts: On-Demand Expert Allocation for Multimodal Language Models with Mixture of Expert. Research goal: Does AnyExperts' dynamic routing strategy exhibit robustness to distribution shift across different multimodal benchmarks (e.g., from GQA to NLVR2) in terms of accuracy consistency and expert utilization variance?.

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

13 papers retrieved. 7 claims extracted, 7 verified. Tribunal: 8.0/10 → 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
AnyExperts is a novel on-demand, budget-aware dynamic routing framework that allocates a variable total number of expert	✓	0.39
The total slots per token are constrained within a fixed range	✓	0.22
Each slot is filled by either a real expert or a virtual expert, with the virtual share capped at a small maximum (e.g.,	✓	0.31
On general image/video tasks, AnyExperts achieves comparable accuracy with 40% fewer real expert activations	✓	0.27
On text-dense tasks (OCR and NLP), AnyExperts maintains performance while reducing real expert usage by 10%	✓	0.28
Existing multimodal MoE approaches rely on rigid routing strategies typically activating a fixed number of experts per t	✓	0.27
Multimodal Mixture-of-Experts (MoE) models offer a promising path toward scalable and efficient large vision-language sy	✓	0.28

References

- <https://www.semanticscholar.org/paper/96b7423354f33c6991b33897119b68690af79229>
- <http://arxiv.org/abs/2604.00086v1>
- <http://arxiv.org/abs/2511.18314v1>