

SOVEREIGN: Can expert specialization patterns in SMOES be transferred across different vision-language tasks (e.g., capti

SOVEREIGN Research Kernel
Autonomous draft — Owner review required before publication

May 28, 2026

Abstract

A pivotal advancement in the progress of large language models (LLMs) is the emergence of the Mixture-of-Experts (MoE) LLMs. Compared to traditional LLMs, MoE LLMs can achieve higher performance with fewer parameters, but it is still hard to deploy them due to their immense parameter sizes. Different from previous weight pruning methods that rely on specifically designed hardware, this paper mainly aims to enhance the deployment efficiency of MoE LLMs by introducing plug-and-play expert-level sparsification techniques. Specifically, we propose, for the first time to our best knowledge, post-tr

1 Introduction

Analysis of: Not All Experts are Equal: Efficient Expert Pruning and Skipping for Mixture-of-Experts Large Language Models. Research goal: Can expert specialization patterns in SMOES be transferred across different vision-language tasks (e.g., captioning, visual grounding) while maintaining performance improvements under equal parameter count constraints?.

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

12 papers retrieved. 9 claims extracted, 1 verified. Tribunal: 3.0/10 → 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
MoE LLMs achieve a reduction in on-the-fly (active) parameters by choosing only top-k experts for the inference of each	×	0.13
Loading the Mixtral 8x7B model in bf16 format requires at least two A100-80G GPUs.	×	0.02
In the Mixtral 8x7B model, the eight experts constitute around 96% (45B out of 47B) of the total number of parameters.	×	0.03
Not all experts are equal in the MoE model.	×	0.13
Unlike existing post-training weight pruning schemes for LLMs, which primarily target unstructured sparsity and N:M semi	×	0.11
Fine-grained weight pruning techniques face challenges in plug-and-play deployment due to the necessity for specific har	×	0.11
Our proposed method significantly reduces memory usage for deploying MoE LLMs and enhances their inference speed.	×	0.11
We examine expert-level pruning for both task-agnostic and task-specific models.	✓	0.18
In the Mixtral 8x7B model, each token x in the input sequence is routed to the top-2 experts based on the routing weight	×	0.09

References

- <http://arxiv.org/abs/2410.17954v2>
- <http://arxiv.org/abs/2402.14800v2>

- <http://arxiv.org/abs/2210.09263v1>