

# SOVEREIGN: What is the computational overhead of implementing expert bridging versus full fine-tuning in terms of inferen

SOVEREIGN Research Kernel

Autonomous draft — Owner review required before publication

May 28, 2026

## Abstract

Fine-tuning Large Language Models (LLMs) is a common practice to adapt pre-trained models for specific applications. While methods like LoRA have effectively addressed GPU memory constraints during fine-tuning, their performance often falls short, especially in multi-task scenarios. In contrast, Mixture-of-Expert (MoE) models, such as Mixtral 8x7B, demonstrate remarkable performance in multi-task learning scenarios while maintaining a reduced parameter count. However, the resource requirements of these MoEs remain challenging, particularly for consumer-grade GPUs with less than 24GB memory. To

## 1 Introduction

Analysis of: MixLoRA: Enhancing Large Language Models Fine-Tuning with LoRA based Mixture of Experts. Research goal: What is the computational overhead of implementing expert bridging versus full fine-tuning in terms of inference latency and memory usage across different model scales from 110M to 175B parameters?.

## 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

19 papers retrieved. 10 claims extracted, 10 verified. Tribunal: 7.7/10 → 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
MixLoRA improves about 9% accuracy compared to state-of-the-art PEFT methods in multi-task learning scenarios.	✓	0.31
MixLoRA reduces GPU memory consumption by 40% during training and inference.	✓	0.19
MixLoRA reduces token computation latency by 30% during training and inference.	✓	0.16
MixLoRA inserts multiple LoRA-based experts within the feed-forward network block of a frozen pre-trained dense model.	✓	0.34
MixLoRA employs a commonly used top-k router.	✓	0.16
MixLoRA utilizes independent attention-layer LoRA adapters to enhance model performance.	✓	0.17
An auxiliary load balance loss is employed to address the imbalance problem of the router.	✓	0.24
Mixtral 8x7B demonstrates remarkable performance in multi-task learning scenarios while maintaining a reduced parameter	✓	0.27
LoRA methods often fall short in performance, especially in multi-task scenarios.	✓	0.17
MoE models like Mixtral 8x7B have resource requirements that remain challenging for consumer-grade GPUs with less than 2	✓	0.24

### References

- <https://www.semanticscholar.org/paper/f0891e0685b1cf25fe184c8a52646099113d9ae0>

- <https://www.semanticscholar.org/paper/dcedb2b0f21d5731144d6475363b37deaf634ce0>
- <https://www.semanticscholar.org/paper/ebcf108f8bc42140721ff02b6727b0a291362957>