

Dynamic Expert Routing in MoE Models Enhances Cross-Domain Generalization on GLUE

Assignee Research

May 30, 2026

Abstract

This report synthesises findings from 16 peer-reviewed papers addressing the following research question: How do MoE-based language models trained with dynamic expert routing perform on cross-domain generalization tasks (measured by GLUE benchmark accuracy) compared to fixed-capacity MoE models and dense. Recent large language models such as Gemini-1.5, DeepSeek-V3, and Llama-4 increasingly adopt Mixture-of-Experts (MoE) architectures, which offer strong efficiency-performance trade-offs by activating only a fraction of the model per token. Yet academic researchers still lack a. 14 claims were extracted from source literature; 1 was independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.8/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: FLAME-MoE: A Transparent End-to-End Research Platform for Mixture-of-Experts Language Models. Research question: How do MoE-based language models trained with dynamic expert routing perform on cross-domain generalization tasks (measured by GLUE benchmark accuracy) compared to fixed-capacity MoE models and dense transformers of equivalent parameter efficiency?.

2 Methodology

Systematic literature search across multiple databases yielded 16 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 5.8/10.

3 Results

16 papers retrieved. 14 claims extracted; 1 independently verified. Quality review score: 5.8/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

Claim	Verified	Confidence
FLAME-MoE significantly outperforms dense counterparts with the same pretraining FLOPs on almost every task.	×	0.05
FLAME-MoE achieves more than 3 points of average accuracy improvements over dense baselines under both 8.0e19 and 2.4e20	×	0.09
FLAME-MoE-290M-1.3B achieves an average downstream accuracy of 0.4154 at 2.0e19 FLOPs.	×	0.05
Dense-411M achieves an average downstream accuracy of 0.4050 at 6.0e18 FLOPs.	×	0.03
FLAME-MoE-1.7B-10.3B achieves an average downstream accuracy of 0.48 at 2.4e20 FLOPs.	×	0.06
Dense-1.4B achieves an average downstream accuracy of 0.44 at 2.4e20 FLOPs.	×	0.03
FLAME-MoE-1.7B-10.3B achieves a throughput of approximately 350 TFLOPS/s/GPU at 100 training steps with PP=1 and EP=8.	×	0.04
Dense-1.4B achieves a throughput of approximately 140 TFLOPS/s/GPU at 100 training steps.	×	0.02
FLAME-MoE-1.7B-10.3B achieves an elapsed time per step of approximately 20 seconds at 100 training steps with PP=1 and E	×	0.04
Dense-1.4B achieves an elapsed time per step of approximately 60 seconds at 100 training steps.	×	0.02
FLAME-MoE includes seven decoder-only MoE models ranging from 38M to 1.7B active parameters.	✓	0.20
FLAME-MoE models use 64 experts per layer, top-8 gating, and two shared experts.	×	0.15
FLAME-MoE is the only MoE platform offering full openness—code, data, checkpoints, routing logs, and evaluation results—	×	0.12
Empirical evaluations on 6 downstream tasks show that FLAME-MoE consistently outperforms dense counterparts trained unde	×	0.07

References

- <http://arxiv.org/abs/2505.20225v1>
- <http://arxiv.org/abs/2604.23996v1>
- <http://arxiv.org/abs/2402.14800v2>