

Sparse Mixture-of-Experts vs. Dense Transformers in Mathematical Reasoning Benchmarks

Assignee Research

June 6, 2026

Abstract

This report synthesises findings from 14 peer-reviewed papers addressing the following research question: How do sparse mixture-of-experts models compare to dense transformers on mathematical reasoning v17. 18 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.2/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Task-Conditioned Routing Signatures in Sparse Mixture-of-Experts Transformers. Research question: How do sparse mixture-of-experts models compare to dense transformers on mathematical reasoning v17.

2 Methodology

Systematic literature search across multiple databases yielded 14 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.2/10.

3 Results

14 papers retrieved. 18 claims extracted; 0 independently verified. Quality review score: 4.2/10.

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
The OLMoE-1B-7B-0125-Instruct model contains 16 MoE layers.	×	0.13
The OLMoE-1B-7B-0125-Instruct model has 64 experts per layer.	×	0.13
The model uses top-k routing with $k = 8$.	×	0.04
Only 8 of 64 experts are active per token in each MoE layer, corresponding to a sparsity level of 12.5%.	×	0.05
The prompt dataset consists of 80 prompts across four categories.	×	0.04
The Code category contains 20 prompts consisting of programming tasks and algorithmic prompts.	×	0.03
The Math category contains 20 prompts consisting of mathematical and symbolic reasoning prompts.	×	0.07
The Story category contains 20 prompts consisting of creative writing and narrative prompts.	×	0.03
The Factual category contains 20 prompts consisting of knowledge retrieval and question-answering prompts.	×	0.05
Each prompt generated 32 tokens during inference.	×	0.04
Within-category routing similarities lie between 0.83 and 0.85.	×	0.05
Cross-category routing similarities typically lie between 0.58 and 0.64.	×	0.06
The empirical routing similarity ordering is Within > LoadBalance > Across.	×	0.09
Task separation measured by Cohen’s d is weakest in early layers and strongest in deeper layers.	×	0.09
The layer-wise task signal (Cohen’s d) peaks around layer 13.	×	0.09
PCA projection of routing signatures shows distinct clusters for code, math, story, and factual prompts.	×	0.08
Story prompts occupy a clearly separated region in the PCA projection.	×	0.01
Code and math prompts form different but partially adjacent clusters in the PCA projection.	×	0.03

References

- <http://arxiv.org/abs/2509.25160v1>
- <http://arxiv.org/abs/2603.11114v1>
- <http://arxiv.org/abs/2402.14800v2>