

# DeepSeek-V3 Shared Expert Mechanism and Its Latency-Throughput Trade-offs on Consumer GPUs

Assignee Research

June 6, 2026

## Abstract

This report synthesises findings from 10 peer-reviewed papers addressing the following research question: How does the shared expert mechanism in DeepSeek-V3 impact latency and token throughput on consumer-grade GPUs compared to standard sparse MoE baselines like Mixtral 8x7B. 11 claims were extracted from source literature; 1 was independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.5/10. 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 does the shared expert mechanism in DeepSeek-V3 impact latency and token throughput on consumer-grade GPUs compared to standard sparse MoE baselines like Mixtral 8x7B?.

## 2 Methodology

Systematic literature search across multiple databases yielded 10 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.5/10.

## 3 Results

10 papers retrieved. 11 claims extracted; 1 independently verified. Quality review score: 4.5/10.

## 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 matches or even outperforms dense models trained with 2x FLOPs (e.g., in 400M-4x).	×	0.06
FLAME-MoE substantially improves pretraining efficiency, achieving a better speed-quality frontier.	×	0.04
Increasing EP generally improves utilization and reduces latency, while deeper pipeline parallelism (e.g., PP=2) can fur	×	0.03
The best-performing configuration for training FLAME-MoE models is PP=1 and EP=8.	×	0.07
FLAME-MoE models demonstrate great utilization under EP=8.	×	0.05
The overall FLOPs throughput of MoE models still lags behind dense models.	×	0.04
FLAME-MoE includes seven decoder-only MoE models (38M–1.7B active parameters), each with 64 experts per layer, top-8 gat	✓	0.21
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/2602.00879v1>
- <http://arxiv.org/abs/2410.07348v1>
- <http://arxiv.org/abs/2505.20225v1>