

DeepSeek-V3 Cross-Domain Finetuning Trade-offs on GPQA Diamond Inference Efficiency

Assignee Research

May 30, 2026

Abstract

This report synthesises findings from 10 peer-reviewed papers addressing the following research question: What is the inference efficiency trade-off when applying cross-domain finetuning to DeepSeek-V3 on GPQA Diamond tasks. We present DeepSeek-V3, a strong Mixture-of-Experts (MoE) language model with 671B total parameters with 37B activated for each token. To achieve efficient inference and cost-effective training, DeepSeek-V3 adopts Multi-head Latent Attention (MLA) and DeepSeekMoE architectures. 9 claims were extracted from source literature; 9 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 9.0/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: DeepSeek-V3 Technical Report. Research question: What is the inference efficiency trade-off when applying cross-domain finetuning to DeepSeek-V3 on GPQA Diamond tasks?.

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 9.0/10.

3 Results

10 papers retrieved. 9 claims extracted; 9 independently verified. Quality review score: 9.0/10.

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
DeepSeek-V3 is a Mixture-of-Experts (MoE) language model with 671B total parameters and 37B activated for each token.	✓	0.30
DeepSeek-V3 adopts Multi-head Latent Attention (MLA) and DeepSeekMoE architectures for efficient inference and cost-efficiency.	✓	0.34
DeepSeek-V3 pioneers an auxiliary-loss-free strategy for load balancing and sets a multi-token prediction training objective.	✓	0.40
DeepSeek-V3 is pre-trained on 14.8 trillion diverse and high-quality tokens.	✓	0.21
DeepSeek-V3 undergoes Supervised Fine-Tuning and Reinforcement Learning stages to fully harness its capabilities.	✓	0.24
DeepSeek-V3 outperforms other open-source models and achieves performance comparable to leading closed-source models.	✓	0.31
DeepSeek-V3 requires only 2.788M H800 GPU hours for its full training.	✓	0.27
The training process of DeepSeek-V3 is remarkably stable with no irrecoverable loss spikes or rollbacks.	✓	0.25
The model checkpoints for DeepSeek-V3 are available at https://github.com/deepseek-ai/DeepSeek-V3 .	✓	0.26

References

- <https://doi.org/10.1038/s41586-025-09422-z>
- <https://doi.org/10.48550/arxiv.2505.09388>
- <https://doi.org/10.48550/arxiv.2412.19437>