

Adaptive Sampling Strategies Reduce Latency and Memory in Federated LLM Deployment

Assignee Research

May 31, 2026

Abstract

This report synthesises findings from 4 peer-reviewed papers addressing the following research question: What is the impact of adaptive sampling strategies on the inference latency and memory footprint of deployed personalized LLMs compared to random sampling in bandwidth-constrained federated networks. Abstract This paper critically examines model compression techniques within the machine learning (ML) domain, emphasizing their role in enhancing model efficiency for deployment in resource-constrained environments, such as mobile devices, edge computing, and Internet of Things. 6 claims were extracted from source literature; 6 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 7.2/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: A comprehensive review of model compression techniques in machine learning. Research question: What is the impact of adaptive sampling strategies on the inference latency and memory footprint of deployed personalized LLMs compared to random sampling in bandwidth-constrained federated networks?.

2 Methodology

Systematic literature search across multiple databases yielded 4 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.2/10.

3 Results

4 papers retrieved. 6 claims extracted; 6 independently verified. Quality review score: 7.2/10.

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
Model compression techniques enhance model efficiency for deployment in resource-constrained environments such as mobile	✓	0.32
Machine learning models are growing increasingly complex and data-intensive.	✓	0.16
The demand for computational resources and memory has surged as machine learning models grow more complex.	✓	0.23
Limited hardware capabilities present significant challenges for deploying AI systems in real-world applications.	✓	0.17
Model compression techniques are essential for utilizing ML models across various domains while maintaining high perform	✓	0.30
Hybrid methods that combine multiple compression techniques promise to deliver superior performance and efficiency compa	✓	0.23

References

- <https://doi.org/10.1007/s10489-024-05747-w>
- <https://doi.org/10.48550/arxiv.2401.08092>
- <https://doi.org/10.48550/arxiv.2303.16129>