

# Quantization-Aware Training Effects on Pruned Transformer Reasoning Under Latency Constraints

Assignee Research

May 30, 2026

## Abstract

This report synthesises findings from 16 peer-reviewed papers addressing the following research question: What is the impact of quantization-aware training on the reasoning capabilities of pruned Transformers compared to full-precision models when measured by MBPP pass@k scores under latency constraints. Large pre-trained transformers are show-stealer in modern-day deep learning, and it becomes crucial to comprehend the parsimonious patterns that exist within them as they grow in scale. With exploding parameter counts, Lottery Ticket Hypothesis (LTH) and its variants, have lost. 16 claims were extracted from source literature; 1 was independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.8/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: The Emergence of Essential Sparsity in Large Pre-trained Models: The Weights that Matter. Research question: What is the impact of quantization-aware training on the reasoning capabilities of pruned Transformers compared to full-precision models when measured by MBPP pass@k scores under latency constraints?.

## 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 3.8/10.

### **3 Results**

16 papers retrieved. 16 claims extracted; 1 independently verified. Quality review score: 3.8/10.

### **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
The paper evaluates the emergence of essential sparsity in large pre-trained models and identifies which weights matter	✓	0.28
Experiments were conducted on multiple NLP benchmarks including MNLI, QNLI, QQP, RTE, SST-2, and SQuAD v1.1.	×	0.02
Computer vision experiments were conducted on CIFAR-10, CIFAR-100, Fashion-MNIST, and Tiny-ImageNet datasets.	×	0.03
The MNLI dataset used 392,704 training examples over 3 epochs with batch size 32 and learning rate $2e-5$ .	×	0.02
The QNLI dataset used 104,768 training examples over 4 epochs with batch size 32 and learning rate $2e-5$ .	×	0.02
The QQP dataset used 363,872 training examples over 3 epochs with batch size 32 and learning rate $2e-5$ .	×	0.02
The RTE dataset used 2,496 training examples over 5 epochs with batch size 32 and learning rate $2e-5$ .	×	0.02
The SST-2 dataset used 67,360 training examples over 5 epochs with batch size 32 and learning rate $2e-5$ .	×	0.02
The SQuAD v1.1 dataset used 88,656 training examples over 3 epochs with batch size 16 and learning rate $3e-5$ .	×	0.01
The CIFAR-10 dataset used 45,000 training examples over 8 epochs with batch size 64 and learning rate $2e-5$ .	×	0.02
The CIFAR-100 dataset used 45,000 training examples over 8 epochs with batch size 64 and learning rate $2e-5$ .	×	0.02
The Fashion-MNIST dataset used 55,000 training examples over 8 epochs with batch size 64 and learning rate $2e-5$ .	×	0.01
The Tiny-ImageNet dataset used 90,000 training examples over 5 epochs with batch size 64 and learning rate $2e-5$ .	×	0.01
AdamW optimizer with decay parameter $\alpha = 1 \times 10^{-8}$ was used for NLP tasks.	×	0.03
AdamW optimizer with decay parameter $\alpha = 2 \times 10^{-8}$ was used for computer vision tasks.	×	0.05
Evaluation metrics include Matched Accuracy for MNLI, Accuracy for QNLI and RTE, F1-score for QQP, and Top-1 Accuracy fo	×	0.04

## References

- <http://arxiv.org/abs/2412.15846v1>
- <http://arxiv.org/abs/2306.03805v2>
- <http://arxiv.org/abs/2210.06313v2>