

Counterfactual Augmentation in LLM Reasoning Under Temporal Shift

Assignee Research

June 11, 2026

Abstract

We give simpler, sparser, and faster algorithms for differentially private fine-tuning of large-scale pre-trained language models, which achieve the state-of-the-art privacy versus utility tradeoffs on many standard NLP tasks. We propose a meta-framework for this problem, inspired by the recent success of highly parameter-efficient methods for fine-tuning. Our experiments show that differentially private adaptations of these approaches outperform previous private algorithms in three important dimensions: utility, privacy, and the computational and memory cost of private training. On many commo

1 Introduction

This paper examines: Differentially Private Fine-tuning of Language Models. Research question: How does counterfactual data augmentation affect LLM reasoning accuracy on the BigBench Hard suite under temporal distribution shift compared to standard fine-tuning?.

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 7.4/10.

3 Results

14 papers retrieved. 17 claims extracted; 11 independently verified. Quality review score: 7.4/10.

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
Privately fine-tuning RoBERTa-Large on the MNLI dataset with a privacy budget of ($\epsilon = 6.7$, $\delta = 1e-6$) achieves an accuracy	✓	0.27
Non-private fine-tuning of RoBERTa-Large on the MNLI dataset achieves an accuracy of 90.2%.	✓	0.24
GPT-3 achieves an accuracy of 91.7% on the MNLI dataset.	×	0.12
Privately fine-tuning GPT-2-Large on the E2E dataset with a privacy budget of ($\epsilon = 6.0$, $\delta = 1e-5$) achieves a ROUGE-L score	✓	0.23
Non-private fine-tuning of GPT-2-Large on the E2E dataset achieves a ROUGE-L score of 72.0.	✓	0.19
On the MNLI dataset with a privacy budget of ($\epsilon = 6.7$, $\delta = 1e-6$), RoBERTa-Base achieves an accuracy of 83.5% compared to	✓	0.25
On the MNLI dataset with a privacy budget of ($\epsilon = 6.7$, $\delta = 1e-6$), RoBERTa-Large achieves an accuracy of 87.8% compared to	✓	0.28
Privately fine-tuning GPT-2-Medium on DART with ($\epsilon = 6.8$, $\delta = 1e-5$) achieves a BLEU score of 42.0.	✓	0.20
Non-private fine-tuning of GPT-2-Medium on DART achieves a BLEU score of 47.1.	✓	0.18
Privately fine-tuning GPT-2-XL on DART with ($\epsilon = 6.8$, $\delta = 1e-5$) achieves a BLEU score of 43.8.	✓	0.20
Non-private fine-tuning of GPT-2-XL on DART achieves a BLEU score of 48.1.	✓	0.17
Full fine-tuning using DPSGD requires 27.9 GB of memory and has a speed of 715 seconds per epoch.	×	0.06
The RGP method requires 9.1 GB of memory and has a speed of 296 seconds per epoch.	×	0.03
The DP LoRA method requires 6.1 GB of memory and has a speed of 271 seconds per epoch.	×	0.04
Using LoRA DP ($r=16$) on RoBERTa-Large results in training 0.94% of the total parameters while achieving an average accuracy	×	0.14
Using Compacter DP ($r=96$, $n=8$) results in training 0.055% of the total parameters.	×	0.06
Prior work has documented that privacy requirements can increase training time by up to two orders of magnitude.	✓	0.17

References

- <http://arxiv.org/abs/2508.04848v1>
- <http://arxiv.org/abs/2110.06500v2>
- <http://arxiv.org/abs/2502.12896v5>