

Parameter-Efficient Fine-Tuning of Small Language Models: A Cross-Domain Evaluation Using LoRA

Aryan Pawar
Mahindra University
aryanpawar1515@mahindra.edu.in

November 22, 2025

Abstract

The increasing computational demands of fine-tuning large language models have motivated the development of parameter-efficient fine-tuning (PEFT) methods. This work presents a comprehensive empirical evaluation of Low-Rank Adaptation (LoRA) applied to Pythia-410M, a small language model (SLM) with 410 million parameters, across five diverse text domains. We systematically evaluate model performance on TinyStories, SQuAD, Yelp Reviews, WikiText-2, and AG News datasets, achieving validation perplexities ranging from 6.46 to 30.55. Our results demonstrate that LoRA enables effective domain adaptation while updating only 0.39% of model parameters (1.57M trainable parameters), reducing training time to approximately 12 minutes per dataset on NVIDIA T4 GPU. We observe significant performance variation across domains, with narrative text (TinyStories, PPL=6.46) showing superior adaptation compared to news classification tasks (AG News, PPL=30.55). All experiments are tracked using Weights & Biases, providing reproducible benchmarks for future PEFT research on small-scale language models.

Keywords: Parameter-Efficient Fine-Tuning, LoRA, Small Language Models, Transfer Learning, Domain Adaptation, Pythia

1 Introduction

1.1 Motivation

The transformer architecture [1] has revolutionized natural language processing, enabling models like GPT-2 [2] and GPT-3 [3] to achieve remarkable performance across diverse tasks. However, the standard approach of full fine-tuning—updating all model parameters—becomes computationally prohibitive as models scale to billions of parameters. For a 175B parameter model like GPT-3, full fine-tuning requires storing and updating all parameters, consuming substantial memory and computational resources.

Parameter-Efficient Fine-Tuning (PEFT) methods address this challenge by updating only a small fraction of model parameters while maintaining competitive performance. Among these approaches, Low-Rank Adaptation (LoRA) [4] has emerged as particularly effective, introducing trainable low-rank decomposition matrices into transformer layers while keeping the original model frozen. This enables fine-tuning with a fraction of the memory and computational cost of full fine-tuning.

1.2 Research Questions

This work investigates the following research questions:

1. **RQ1:** How effectively does LoRA enable domain adaptation for small language models across diverse text types?
2. **RQ2:** What is the computational efficiency of LoRA fine-tuning in terms of training time, memory usage, and parameter count?
3. **RQ3:** How does model performance vary across different text domains, and what factors influence this variation?
4. **RQ4:** Can LoRA-based fine-tuning achieve competitive performance with minimal parameter updates (<1% trainable parameters)?

1.3 Contributions

Our contributions are as follows:

- **Comprehensive Evaluation:** We conduct systematic experiments across five diverse datasets spanning narratives, question-answering, reviews, encyclopedic text, and news classification.
- **Efficiency Analysis:** We demonstrate that LoRA enables effective fine-tuning with only 0.39% trainable parameters and \sim 12 minute training times per dataset.
- **Cross-Domain Insights:** We provide detailed analysis of performance patterns across domains, identifying text characteristics that correlate with adaptation success.
- **Reproducible Framework:** We release a complete implementation with experiment tracking, enabling reproduction and extension of our results.

2 Background and Related Work

2.1 Transformer Language Models

The transformer architecture [1] introduced self-attention mechanisms that enable modeling long-range dependencies in sequences. Modern language models like GPT-2 [2] and GPT-3 [3] demonstrate that scaling model size and training data leads to emergent capabilities in few-shot and zero-shot learning. However, this scaling comes with increased computational requirements for both training and fine-tuning.

Small Language Models (SLMs) in the 100M-1B parameter range, such as Pythia [10], offer a practical middle ground. While less capable than larger models, SLMs can be trained and deployed with modest computational resources, making them accessible for research and production use.

2.2 Parameter-Efficient Fine-Tuning

Full fine-tuning updates all model parameters, requiring storage of gradients and optimizer states for every parameter. For large models, this becomes impractical. PEFT methods address this by introducing small trainable modules while keeping most of the model frozen.

Several PEFT approaches have been proposed:

Adapter Layers [8]: Insert small feedforward networks between transformer layers. Adapter-Fusion [6] extends this by learning to compose multiple adapters.

Prefix Tuning [9]: Prepends trainable continuous prompts to input sequences, optimizing these prompts while keeping the model frozen.

BitFit [7]: Updates only bias terms in the model, achieving surprising effectiveness with minimal parameter updates.

LoRA [4]: Introduces low-rank decomposition matrices into attention layers, enabling efficient adaptation with competitive performance.

2.3 Low-Rank Adaptation (LoRA)

LoRA [4] hypothesizes that weight updates during fine-tuning have low intrinsic dimensionality. For a pre-trained weight matrix $W_0 \in \mathbb{R}^{d \times k}$, LoRA represents the update as:

$$W = W_0 + \Delta W = W_0 + BA \quad (1)$$

where $B \in \mathbb{R}^{d \times r}$ and $A \in \mathbb{R}^{r \times k}$ with rank $r \ll \min(d, k)$. During training, W_0 remains frozen while A and B are updated. The number of trainable parameters is reduced from $d \times k$ to $r \times (d + k)$.

LoRA offers several advantages:

- **Memory Efficiency:** Only stores gradients for low-rank matrices
- **No Inference Latency:** Merged weights ($W_0 + BA$) at deployment
- **Modular:** Multiple LoRA adapters can be trained for different tasks
- **Competitive Performance:** Achieves results comparable to full fine-tuning

3 Methodology

3.1 Model Architecture

We use **Pythia-410M** [10], a decoder-only transformer language model with the following specifications:

Parameter	Value
Total Parameters	406,851,584
Hidden Dimension	1,024
Number of Layers	24
Attention Heads	16
Vocabulary Size	50,277
Context Length	2,048 tokens

Table 1: Pythia-410M Model Specifications

Pythia is trained on The Pile [11], a diverse 800GB corpus, and provides checkpoints at multiple training stages, enabling controlled experimentation.

3.2 LoRA Configuration

We apply LoRA to the query-key-value projection matrices in all attention layers. Our configuration:

The scaling factor $\alpha/r = 2$ controls the magnitude of LoRA updates. We freeze all original model parameters and update only the low-rank matrices A and B .

3.3 Training Configuration

All experiments use consistent training hyperparameters to ensure fair comparison across datasets:

We use gradient checkpointing to reduce memory consumption and enable training on a single NVIDIA T4 GPU with 16GB VRAM.

Hyperparameter	Value
LoRA Rank (r)	16
LoRA Alpha (α)	32
LoRA Dropout	0.1
Target Modules	<code>query_key_value</code>
Bias	None
Trainable Parameters	1,572,864
Total Parameters	406,851,584
Trainable %	0.39%

Table 2: LoRA Hyperparameters and Parameter Counts

Hyperparameter	Value
Optimizer	AdamW
Learning Rate	1×10^{-4}
LR Scheduler	Cosine
Warmup Steps	100
Weight Decay	0.01
Batch Size	8
Gradient Accumulation	1
Max Gradient Norm	1.0
Epochs	1
Max Sequence Length	512 tokens
Precision	FP32
Random Seed	42

Table 3: Training Hyperparameters

3.4 Datasets

We evaluate on five datasets covering diverse domains:

TinyStories [12]: A dataset of simple, child-appropriate narratives generated synthetically. Contains 2.1M training examples. We sample 2,000 examples for efficiency.

SQuAD [13]: Stanford Question Answering Dataset containing Wikipedia passages for reading comprehension. We use context passages as training data (1,000 samples).

Yelp Reviews: Restaurant reviews from the Yelp dataset, spanning ratings from 1-5 stars (1,000 samples).

WikiText-2 [14]: A collection of verified Wikipedia articles, representing formal encyclopedic writing (1,000 samples).

AG News [15]: News article classification dataset with four categories: World, Sports, Business, Sci/Tech (2,000 samples).

All datasets are tokenized using the Pythia tokenizer with a maximum length of 512 tokens. We split each dataset 90/10 for training and validation.

3.5 Evaluation Metrics

We report three primary metrics:

Dataset	Domain	Samples	Avg Length
TinyStories	Narratives	2,000	Short
SQuAD	QA Context	1,000	Medium
Yelp Reviews	Reviews	1,000	Short
WikiText-2	Encyclopedia	1,000	Long
AG News	News	2,000	Short

Table 4: Dataset Characteristics

Validation Loss: Cross-entropy loss on held-out validation data, computed as:

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \log P(x_i | x_{<i}) \quad (2)$$

Perplexity (PPL): Exponential of validation loss, measuring model uncertainty:

$$PPL = \exp(\mathcal{L}) \quad (3)$$

Lower perplexity indicates better language modeling performance.

Bits Per Token (BPT): Information-theoretic measure of compression:

$$BPT = \frac{\mathcal{L}}{\log 2} \quad (4)$$

3.6 Implementation Details

Our implementation uses:

- **Framework:** PyTorch 2.0, Hugging Face Transformers, PEFT library
- **Hardware:** Single NVIDIA T4 GPU (16GB VRAM), 2 CPU cores
- **Experiment Tracking:** Weights & Biases for logging metrics, system stats
- **Training Time:** ~12 minutes per dataset, ~1 hour total

Code is available at: <https://github.com/aryanpawar1234/slm-lora-benchmark>

4 Experimental Results

4.1 Overall Performance

Table 5 presents the primary results across all five datasets. We observe substantial variation in performance, with perplexity ranging from 6.46 (TinyStories) to 30.55 (AG News).

4.2 Domain-Specific Analysis

Narratives (TinyStories): Achieves the lowest perplexity (6.46), indicating strong adaptation. The simple, repetitive structure of children’s stories aligns well with the model’s pre-training on natural language.

Question-Answering (SQuAD): Moderate perplexity (14.68) on Wikipedia contexts. The formal, informative writing style presents moderate difficulty.

Reviews (Yelp): Higher perplexity (25.40) reflects the informal, varied nature of user-generated content with diverse vocabulary and colloquialisms.

Dataset	Domain	Loss ↓	PPL ↓	BPT ↓	Time (min)
TinyStories	Narratives	1.87	6.46	2.69	12
SQuAD	QA Context	2.69	14.68	3.88	11
Yelp Reviews	Reviews	3.23	25.40	4.67	11
WikiText-2	Encyclopedia	3.41	30.28	4.92	11
AG News	News	3.42	30.55	4.93	11
Mean	—	2.92	21.47	4.22	11.2
Std Dev	—	0.68	10.67	0.98	0.4

Table 5: Performance across five datasets. All models trained with identical LoRA configuration ($r=16$, $\alpha=32$). Lower values indicate better performance.

Encyclopedic Text (WikiText-2): High perplexity (30.28) despite Wikipedia being in the pre-training corpus. The specific articles in WikiText-2 may require domain-specific knowledge.

News Classification (AG News): Highest perplexity (30.55), possibly due to the classification-oriented nature of the task and diverse news topics requiring specialized vocabulary.

4.3 Computational Efficiency

Table 6 demonstrates the efficiency gains of LoRA:

Metric	Full Fine-Tuning	LoRA
Trainable Parameters	406,851,584 (100%)	1,572,864 (0.39%)
Training Time (per epoch)	~2 hours*	11-12 min
GPU Memory	~16 GB	~8 GB
Checkpoint Size	1.6 GB	6 MB

Table 6: Efficiency comparison. *Estimated based on typical full fine-tuning times.

Key efficiency observations:

- **Parameter Reduction:** $257\times$ fewer trainable parameters
- **Speed:** $\sim 10\times$ faster training per epoch
- **Memory:** 50% reduction in GPU memory usage
- **Storage:** $266\times$ smaller adapter checkpoints

4.4 Training Dynamics

Figure 1 shows training loss curves for all datasets (placeholder - add W&B screenshots):

[Insert W&B training loss curves here]

5 subplots showing loss vs. steps for each dataset

Figure 1: Training loss curves across datasets. All models converge within 200-230 steps.

All models converge smoothly without signs of overfitting or instability. The cosine learning rate schedule enables stable optimization.

5 Discussion

5.1 Performance Patterns

Our results reveal clear performance patterns related to text characteristics:

Text Simplicity: Simpler, more predictable text (TinyStories) achieves lower perplexity. The repetitive narrative structure and limited vocabulary reduce modeling difficulty.

Domain Specificity: Specialized domains (AG News, WikiText-2) show higher perplexity, suggesting the need for domain-specific knowledge not fully captured by general pre-training.

Formality vs. Informality: Formal text (SQuAD, WikiText-2) shows moderate perplexity, while informal user-generated content (Yelp) presents challenges due to linguistic diversity.

5.2 LoRA Effectiveness

Despite updating only 0.39% of parameters, LoRA achieves effective adaptation across all domains. This validates the low-rank hypothesis: weight updates during fine-tuning lie in a low-dimensional subspace.

The consistent training times (\sim 11-12 minutes) across datasets demonstrate LoRA’s scalability. The small adapter size (6 MB) enables efficient storage and deployment of multiple domain-specific models.

5.3 Comparison to Related Work

While direct comparison is difficult due to different base models and datasets, our results align with findings from the LoRA paper [4], which reported competitive performance with full fine-tuning using similar parameter budgets.

Our work extends prior research by:

- Evaluating small models (410M parameters vs. GPT-3’s 175B)
- Covering diverse text domains in a single study
- Providing reproducible benchmarks with open-source implementation

5.4 Limitations

Several limitations should be noted:

Single Epoch Training: We train for only one epoch due to computational constraints. Additional epochs may improve performance.

Limited Hyperparameter Search: We use a single LoRA rank ($r=16$). Optimal ranks may vary by dataset.

Small Dataset Samples: We use 1,000-2,000 samples per dataset. Larger samples may yield different results.

Single Model: We evaluate only Pythia-410M. Results may differ for other model architectures and sizes.

6 Conclusion and Future Work

6.1 Summary

This work presents a comprehensive evaluation of LoRA-based parameter-efficient fine-tuning for small language models across five diverse text domains. Our key findings:

1. LoRA enables effective domain adaptation with only 0.39% trainable parameters

2. Performance varies significantly across domains (PPL: 6.46-30.55)
3. Training efficiency is excellent: \sim 12 minutes per dataset on T4 GPU
4. Simpler, more structured text (narratives) shows superior adaptation

These results demonstrate that PEFT methods like LoRA make fine-tuning accessible even for resource-constrained settings, enabling practical deployment of domain-adapted language models.

6.2 Future Directions

Several directions for future work emerge:

Hyperparameter Optimization: Systematic search over LoRA rank, learning rate, and other hyperparameters may improve performance.

Multi-Epoch Training: Extending training beyond one epoch with appropriate regularization to prevent overfitting.

Larger Models: Evaluating whether our findings generalize to larger models (1B-7B parameters).

Task-Specific Evaluation: Beyond perplexity, evaluate on downstream tasks (classification, generation quality).

Alternative PEFT Methods: Compare LoRA to adapters, prefix tuning, and prompt tuning.

Mixture of Adapters: Investigate combining multiple domain-specific LoRA adapters for multi-task performance.

Quantization: Combine LoRA with 4-bit or 8-bit quantization for further efficiency gains.

6.3 Broader Impact

By demonstrating effective fine-tuning with minimal resources, this work contributes to democratizing access to language model adaptation. Researchers and practitioners with limited computational budgets can leverage PEFT methods to create domain-specific models, fostering innovation in resource-constrained settings.

Acknowledgments

We thank Weights & Biases for providing experiment tracking infrastructure, and the Hugging Face team for the Transformers and PEFT libraries. This work was conducted using Google Colab’s free GPU resources.

References

- [1] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30.
- [2] Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8), 9.
- [3] Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33, 1877-1901.

- [4] Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., & Chen, W. (2021). LoRA: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- [5] Fedus, W., Zoph, B., & Shazeer, N. (2021). Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *arXiv preprint arXiv:2101.03961*.
- [6] Pfeiffer, J., Kamath, A., Rücklé, A., Cho, K., & Gurevych, I. (2021). AdapterFusion: Non-destructive task composition for transfer learning. *Proceedings of EACL*, 487-503.
- [7] Ben Zaken, E., Ravfogel, S., & Goldberg, Y. (2022). BitFit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. *Proceedings of ACL*, 1-9.
- [8] Houlsby, N., Giurgiu, A., Jastrzebski, S., Morrone, B., De Laroussilhe, Q., Gesmundo, A., Attariyan, M., & Gelly, S. (2019). Parameter-efficient transfer learning for NLP. *Proceedings of ICML*, 2790-2799.
- [9] Li, X. L., & Liang, P. (2021). Prefix-tuning: Optimizing continuous prompts for generation. *Proceedings of ACL-IJCNLP*, 4582-4597.
- [10] Biderman, S., Schoelkopf, H., Anthony, Q., Bradley, H., O'Brien, K., Hallahan, E., ... & Belrose, N. (2023). Pythia: A suite for analyzing large language models across training and scaling. *Proceedings of ICML*, 2397-2430.
- [11] Gao, L., Biderman, S., Black, S., Golding, L., Hoppe, T., Foster, C., ... & Leahy, C. (2020). The Pile: An 800GB dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*.
- [12] Eldan, R., & Li, Y. (2023). TinyStories: How small can language models be and still speak coherent English? *arXiv preprint arXiv:2305.07759*.
- [13] Rajpurkar, P., Zhang, J., Lopyrev, K., & Liang, P. (2016). SQuAD: 100,000+ questions for machine comprehension of text. *Proceedings of EMNLP*, 2383-2392.
- [14] Merity, S., Xiong, C., Bradbury, J., & Socher, R. (2016). Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*.
- [15] Zhang, X., Zhao, J., & LeCun, Y. (2015). Character-level convolutional networks for text classification. *Advances in Neural Information Processing Systems*, 28.

A Hyperparameter Sensitivity

Future work should investigate sensitivity to key hyperparameters:

- LoRA rank: {4, 8, 16, 32, 64}
- Learning rate: {1e-5, 5e-5, 1e-4, 5e-4}
- Batch size: {4, 8, 16, 32}

B Weights & Biases Links

Experiment runs available at:

- TinyStories: <https://wandb.ai/aryanpawar1515-mahindra-university/slm-lora-benchmark/runs/hrr5g7ws>

- AG News: <https://wandb.ai/aryanpawar1515-mahindra-university/slm-lora-benchmark/runs/txpy0k4n>
- Project: <https://wandb.ai/aryanpawar1515-mahindra-university/slm-lora-benchmark>

C Code Availability

Complete implementation available at: <https://github.com/aryanpawar1234/slm-lora-benchmark>