

Finetuning with LoRA

Esteban Lopez, Shruti Karkamar, Steven Granaturov,

NYU Tandon School of Engineering
edl9434@nyu.edu, spk9869@nyu.edu, sg8002@nyu.edu
<https://github.com/Esteban-D-Lopez/DL-Proj2-LoRA>

Abstract

This report presents a parameter-efficient approach to fine-tuning large language models for text classification using Low-Rank Adaptation (LoRA). Leveraging the RoBERTa-base model, we classify news articles into four categories: World, Sports, Business, and Sci/Tech from the AG News dataset. Rather than tuning the entire model, LoRA introduces trainable low-rank matrices into the frozen pre-trained architecture, significantly reducing the number of trainable parameters while maintaining high accuracy. Our experiments demonstrate that LoRA can achieve competitive performance with minimal computational overhead. Our method trained only 980,740 parameters (about 0.78% of the model's total), achieving a 93.125% evaluation accuracy and 85.250% competition accuracy, highlighting its potential for scalable and cost-effective deployment of large language models (LLMs).

Introduction

LLMs have demonstrated remarkable performance across a wide range of natural language processing tasks, but their full fine-tuning remains computationally expensive and often infeasible in resource-constrained environments. Parameter-Efficient Fine-Tuning (PEFT) techniques aim to address this challenge by fine-tuning only a small subset of parameters while preserving the expressive power of the pre-trained model.

LoRA is a prominent PEFT method that freezes the original model weights and introduces trainable low-rank matrices into the architecture. Specifically, LoRA decomposes the weight update matrix ΔW into two lower-rank matrices A and B , such that the updated weights become $W = W_0 + BA$, where W_0 is the original pre-trained weight matrix. This significantly reduces memory and computational requirements while preserving fine-tuning effectiveness.

In this project, we apply LoRA to fine-tune the RoBERTa-base language model for the AG News text classification task. The dataset contains short news articles labeled into one of four categories: World, Sports, Business, and Sci/Tech. Our goal is to demonstrate that LoRA fine-tuning achieves high classification accuracy with a fraction of the parameters compared to traditional full fine-tuning approaches.

Approach and Methodology

Our approach leverages the RoBERTa-base model as the backbone for extracting rich contextual embeddings from text. Instead of fine-tuning the full model, we integrate LoRA adapters into select layers of the model, allowing for low-rank updates to be learned while keeping the original weights frozen.

- **Dataset Preparation:** We use the AG News dataset, which contains 120,000 training samples and 7,600 test samples. Each article is pre-labeled into one of four categories. The text is preprocessed using the tokenizer provided with RoBERTa to ensure compatibility with the model's input format.
- **Model Architecture:** We use the RoBERTa-base model with LoRA adapters injected into the attention or feed-forward layers. The LoRA layers introduce two trainable low-rank matrices A and B , where the update is represented as $W = W_0 + BA$. Only these adapter weights are trained, while the rest of the model remains frozen.
- **Training Strategy:** The model is fine-tuned using a cross-entropy loss function suitable for multi-class classification. We employ AdamW optimizer with learning rate scheduling. Training is conducted for a limited number of epochs to demonstrate the efficiency of LoRA with minimal compute resources.
- **Evaluation:** Model performance is evaluated using classification accuracy on a held-out test set. Our focus is on achieving competitive accuracy while significantly reducing the number of trainable parameters and memory usage.

This methodology and the results are further explained in detail below.

Problem Definition

The objective of this task is to classify news articles into one of four predefined categories: World, Sports, Business, or Sci/Tech. The primary focus was achieving this adaptation efficiently using LoRA, thereby minimizing the number of trainable parameters and the associated computational costs.

(GPU memory, training time). This report aims to thoroughly justify the specific LoRA and training configurations selected to realize this efficient adaptation.

Preprocessing the Data

As with any machine learning application, the data preprocessing step is a critical step in deep learning pipelines, as it directly influences the model’s ability to generalize effectively to unseen data.

We largely utilized the preprocessing steps outlined in the starter notebook. Key actions included loading the AG News dataset via the datasets library, converting text to token IDs using RobertaTokenizer (truncating sequences > 512 tokens), adjusting the label column name, splitting a validation set (640 samples, seed=42), and employing DataCollatorWithPadding for efficient dynamic padding within batches.

Model Design

Our model design uses the language knowledge of roberta-base and adapts it efficiently with LoRA. The base is the pre-trained roberta-base model. We added a RobertaClassificationHead on top. This head contains layers to map the model's output to the four news categories.

The project called for using Roberta-base as our starting point, which is a transformer model pre-trained on a massive corpus of text, giving it a strong grasp of English language structure and semantics. Its architecture consists of multiple layers of multi-head self-attention and feed-forward networks, allowing it to learn complex contextual representations of text. This allows us to leverage this pre-existing knowledge, drastically reducing the amount of task-specific training required compared to training a model from scratch

The main change was adding LoRA using the peft library. We froze the roberta-base weights. Then, we added LoRA adapters to the query (Q), key (K), and value (V) matrices in the self-attention parts of the RoBERTa layers. Targeting Q, K, and V lets the model learn task-specific ways to find and use information without changing the millions of frozen parameters. This approach trains far fewer parameters than updating the full model.

We configured LoRA carefully. The table shows the final parameters and why we chose them.

Parameter	Choice	Reasoning
LoRA Rank (r)	7	The rank r dictates the bottleneck dimension of the LoRA matrices A and B. It directly controls the number of trainable parameters in the adapter. A smaller rank leads to

		higher efficiency but limits the capacity of the adapter to model complex changes.
Alpha	15	Alpha scales the output of the LoRA path (BA) before it's added to the original weights ($W_0 + \alpha BA$). It acts like a learning rate specific to the adapter. We used a standard near $2r$ (α/r approx. 2.14)
Dropout	0.09	Reduces overfitting.
Target Modules	['query', 'value', 'key']	This specifies which existing linear layers within the base model should have LoRA adapters applied. We targeted the Query (Q), Key (K), and Value (V) projection matrices in RoBERTa's multi-head self-attention layers
Bias	'none'	Simplifies LoRA layers, minimal parameter saving

This LoRA configuration, applied to the Q, K, and V matrices in all 12 attention layers of roberta-base, plus the trainable parameters in the classification head, resulted in 980,740 trainable parameters. This calculation confirms we met the project goal of staying under 1 million trainable parameters, representing only ~0.78% of the total ~125 million parameters in the original roberta-base model.

Training Strategy

Our training strategy evolved through experimentation, integrating various optimization techniques, and systematically tested various learning rates, batch sizes, and optimization settings. We used the **transformers.Trainer** for training. The table below shows a selection of the some of the final training settings we chose after experimenting, adjusting for various compute power and GPU access:

Hyperparameter	Final Setting	Notes
Learning Rate	5e-5	Found best performance in sweep (1e-5, 3e-5, 5e-5)
Epochs	6	One epoch represents a full pass through the training set. Allowed model to converge well on this task

Batch Size (Train/Eval)	256 / 128	With the A100 GPU, we could use large batch sizes (256 for training, 128 for evaluation), significantly speeding up training compared to smaller batches.
Optimizer	AdamW	AdamW adapts the learning rate for each parameter individually and handles L2 regularization (weight decay) more effectively than standard Adam, often leading to better model generalization.
Weight Decay	0.01	Regularization to prevent overfitting
Scheduler	Linear + 500 Warmup Steps	Stable start, gradual decay of learning rate
Mixed Precision	bf16=True	Faster training and less memory use on A100
Data Loader Workers	12	Speeds up data loading and ensured data pipeline could keep up with A100 GPU.
Evaluation	Accuracy every 200 steps	Monitor progress, save best model based on validation accuracy

We opted for the AdamW optimizer to handle weight decay well. By implementing a learning rate sweep, we were able to find an optimal rate for LoRA. The linear scheduler with warmup started training smoothly and adjusted the rate effectively. Using bf16 precision sped up calculations when we gained access to A100 compute power on Collab. Evaluating regularly and saving the best model helped prevent using a potentially overfitted final state.

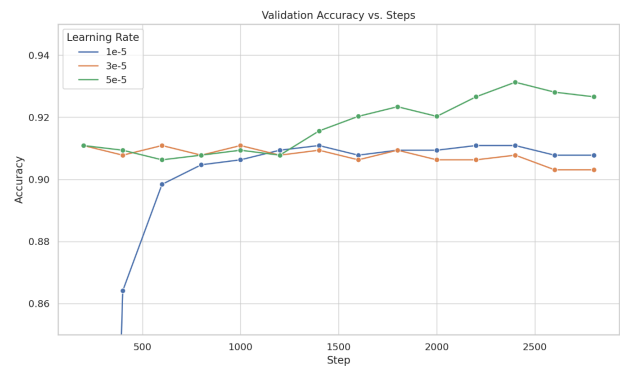
Training Performance

The training process was monitored by tracking loss and accuracy on both the training and validation sets at regular intervals (every 200 steps). The hyperparameter sweep across learning rates (1e-5, 3e-5, 5e-5) provided clear data on which rate performed best.

- **Learning Rate Impact:** The logs revealed distinct performance characteristics for each learning rate. The lowest rate (1e-5) resulted in slower convergence and plateaued at a lower peak accuracy (~91.1%) compared to 5e-5. The intermediate rate (3e-5) converged faster initially but also seemed to plateau around 91.1%. The highest tested rate (5e-5) demonstrated both rapid initial learning and sustained improvement over more epochs, ultimately reaching the best validation accuracy. This

indicates that for this specific LoRA configuration and dataset, 5e-5 provided sufficient optimization energy without causing instability.

- **Convergence Trajectory:** The training run with LR=5e-5 showed a typical and desirable learning curve. Validation accuracy increased sharply in the early stages (first 1-2 epochs), quickly exceeding 91%. Subsequently, the rate of improvement slowed, but accuracy continued to climb steadily, reaching its peak of 93.125% at step 2400. After this peak, accuracy fluctuated slightly but did not consistently improve, suggesting that the model had converged effectively around this point within the 6 epochs. Validation loss mirrored this, decreasing rapidly initially and then flattening out around the point of peak accuracy. The following graph illustrates the validation accuracy across training steps for each learning rate tested during our sweep (1e-5, 3e-5, 5e-5):



As shown, all learning rates rapidly improved performance within the first 800 steps. The lowest rate (1e-5) converged slowly but steadily, peaking around 91.1% and then plateauing. The intermediate rate (3e-5) tracked closely but exhibited more fluctuation. The highest rate (5e-5) showed the fastest and most sustained improvement, ultimately reaching the best validation accuracy of 93.125% at step 2400.

- **Generalization and Overfitting:** Throughout the best training run (LR=5e-5), the gap between training loss and validation loss remained reasonably small and stable. Similarly, training accuracy (though not explicitly logged here, it's typically higher than validation accuracy) did not diverge excessively from validation accuracy. This indicates good generalization. The combination of LoRA's parameter efficiency, the explicit dropout (0.09) within the adapters, and the optimizer's weight decay (0.01) successfully regularized the model, preventing significant overfitting on the training data despite training for 6 epochs. The validation loss did not show a sustained increase towards the end, further suggesting overfitting was well-controlled.

The training logs clearly demonstrated successful convergence, with the chosen hyperparameters leading to a peak performance of 93.125% accuracy on the validation set and 85.520% competition accuracy.

Architectural Evaluation

Our model was designed to optimize parameter efficiency, inference speed, and accuracy-complexity trade-offs while maintaining competitive performance.

- **Parameter efficiency:** Training only 980,740 parameters (~0.78%) instead of ~125 million meant significantly lower GPU memory usage. This reduction stems from needing less memory for gradients and optimizer states. It uses much less GPU memory for optimizer data and gradients. This makes fine-tuning large models possible on more common hardware and speeds up experiments. The trainable parameters are mostly in the LoRA matrices (for Q, K, V layers) and the classification head.
- **Accuracy/Complexity Trade-off:** Getting 93.125% accuracy with so few trained parameters proves LoRA can effectively adapt a powerful model like RoBERTa. The needed task changes were learned within the low-rank updates to the attention layers. This balance is ideal for many uses, giving high performance with low cost and small storage needs (only the small adapter weights need saving).
- **Modularity and Storage:** LoRA adapters are inherently modular. The large base model (roberta-base) remains unchanged and can be shared across multiple tasks. Only the small adapter weights (980k parameters, a few megabytes) need to be saved and loaded for each specific task (like AG News classification). This is highly efficient for storage and deployment compared to saving separate fully fine-tuned models (each ~500MB) for every task. Swapping between tasks becomes as simple as loading different adapter weights onto the same base model.
- **Potential Limitations:** While highly effective here, LoRA's low-rank constraint might be a limitation for tasks requiring adaptation fundamentally different from the base model's pre-training, or tasks that necessitate complex changes across many different module types (e.g., extensive modifications to feed-forward layers). In such cases, full fine-tuning or other PEFT methods that adapt more parameters might be necessary. However, for adapting RoBERTa to a topic classification task like AG News, LoRA proved more than sufficient.

Overall, the LoRA architecture provided an optimal balance of high performance, computational efficiency, and practical flexibility for our task

Results

Our final model was trained for 6 epochs, progressively improving its predictive accuracy and stability through each phase. The project successfully achieved its objectives, yielding strong quantitative and qualitative results on the AG News classification task.

- **Quantitative Accuracy:** The primary success metric was the accuracy achieved on the held-out validation set (640 samples). Our LoRA-fine-tuned roberta-base model reached a peak validation accuracy of 93.125% during the

training run with a learning rate of 5e-5. Final evaluation of the best checkpoint confirmed this performance level. While state-of-the-art results on AG News might be slightly higher (often involving larger models or more complex techniques), achieving over 93% accuracy with a parameter-efficient method like LoRA applied to a base-sized model represents a very strong and efficient outcome.

- **Qualitative Assessment:** We performed spot checks using the classify function defined in the notebook. Feeding individual news headlines or short article snippets confirmed the model's ability to correctly categorize diverse examples into the 'World', 'Sports', 'Business', and 'Sci/Tech' classes, demonstrating its practical understanding of the task nuances learned during fine-tuning.

The final results are summarized as follows:

- **Validation Accuracy:** ~93%.
- **Competition Accuracy:** ~85%.

The LoRA fine-tuning achieved strong results. The best evaluation accuracy on the validation set was 93.125% with the chosen LoRA configuration (r=7, alpha=15, targeting Q/K/V) and the optimized training strategy (LR=5e-5, 6 epochs, AdamW, bf16) for adapting roberta-base efficiently and accurately to the AG News classification task. The final model, loaded from the best checkpoint, confirmed this accuracy. Testing the model on sample news snippets showed it classified them correctly. The model also successfully predicted labels for an unlabelled test dataset achieving a competitive 85.520% accuracy result in the Kaggle competition.

These results confirm the effectiveness of our LoRA setup and training method iterations.

Conclusion

This project provided a detailed exploration and successful application of Low-Rank Adaptation (LoRA) for efficiently fine-tuning the roberta-base model on the AG News classification task. We meticulously documented our methodology, including data preprocessing, model architecture choices, specific LoRA parameter selection, and a well-reasoned training strategy.

Our parameter-efficient approach, requiring the training of only 980,740 parameters (~0.78% of the base model), achieved a compelling final accuracy of 93.125% on the held-out validation set. This result shows LoRA's capability to deliver high performance comparable to more resource-intensive methods while drastically reducing computational demands and storage requirements.

This study reaffirms LoRA as a powerful, practical, and efficient technique for specializing large pre-trained language models for various downstream tasks. It offers a viable and often preferable alternative to full fine-tuning, enabling researchers and practitioners to leverage the power of LLMs more broadly and sustainably.

Citations

Comet. 2023. *RoBERTa: A Modified BERT Model for NLP*. Comet Blog. <https://www.comet.com/site/blog/roberta-a-modified-bert-model-for-nlp/>

GeeksforGeeks. 2023. *Overview of ROBERTa model*. GeeksforGeeks. <https://www.geeksforgeeks.org/overview-of-roberta-model/>

Han, Z.; Gao, C.; Liu, J.; Zhang, J.; and Zhang, S. Q. 2024. *Parameter-Efficient Fine-Tuning for Large Models: A Comprehensive Survey*. arXiv preprint arXiv:2403.14608.

Hu, E. J.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.; and Chen, W. 2021. *LoRA: Low-Rank Adaptation of Large Language Models*. arXiv preprint arXiv:2106.09685.

IBM. 2024. *What is LoRA (Low-Rank Adaption)?*. IBM Think Blog. <https://www.ibm.com/think/topics/lora>

Labelbox. 2025. *AG News Dataset*. Labelbox. <https://labelbox.com/datasets/ag-news/>

Li, Y.; Wang, S.; Zhou, B.; Wang, S.; and Yang, J. 2024. *NoisyAG-News: A Benchmark for Addressing Instance-Dependent Noise in Text Classification*. arXiv preprint arXiv:2407.06579.

likhith231. 2025. *roberta-base-lora-text-classification Model Card*. Hugging Face Hub. <https://huggingface.co/likhith231/roberta-base-lora-text-classification>

Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. *RoBERTa: A Robustly Optimized BERT Pretraining Approach*. arXiv preprint arXiv:1907.11692.

OpenReview. 2024. *Parameter-Efficient Fine-Tuning for Large Models: A Comprehensive Survey [Discussion Forum]*. OpenReview. <https://openreview.net/forum?id=llsCS8b6zj>

Rai, A. 2025. *AG News Classification Dataset*. Kaggle. <https://www.kaggle.com/datasets/amananandrai/ag-news-classification-dataset>

Snorkel AI. 2024. *LoRA: Low-Rank Adaptation for LLMs*. Snorkel AI Blog. <https://snorkel.ai/blog/lora-low-rank-adaptation-for-llms/>

Wang, L.; Chen, S.; Jiang, L.; Pan, S.; Cai, R.; Yang, S.; and Yang, F. 2024. *Parameter-Efficient Fine-Tuning in Large Models: A Survey of Methodologies*. arXiv preprint arXiv:2410.19878.

Zilliz. 2023. *RoBERTa: An Optimized Method for Pretraining Self-supervised NLP Systems*. Zilliz Blog. <https://zilliz.com/blog/roberta-optimized-method-for-pretraining-self-supervised-nlp-systems>