

# Research on LoRA and QLoRA

## Introduction

Large language models (LLMs) have achieved remarkable success across natural language processing tasks. However, their massive size makes full fine-tuning expensive and impractical. Training all parameters of models with billions of weights requires huge compute resources and storage. To address these limitations, researchers proposed parameter-efficient fine-tuning (PEFT) methods. Two of the most impactful techniques are **LoRA (Low-Rank Adaptation)** and **QLoRA (Quantized Low-Rank Adaptation)**. Both reduce memory usage and training cost while preserving performance, making LLM fine-tuning more accessible.

## 1 LoRA (Low-Rank Adaptation)

LoRA, introduced in 2021, is a parameter-efficient fine-tuning method. Instead of updating all the parameters of a pretrained model, LoRA freezes the original model weights and inserts small trainable low-rank matrices into specific layers (mainly attention and feed-forward layers). This enables efficient adaptation of very large models on modest hardware.

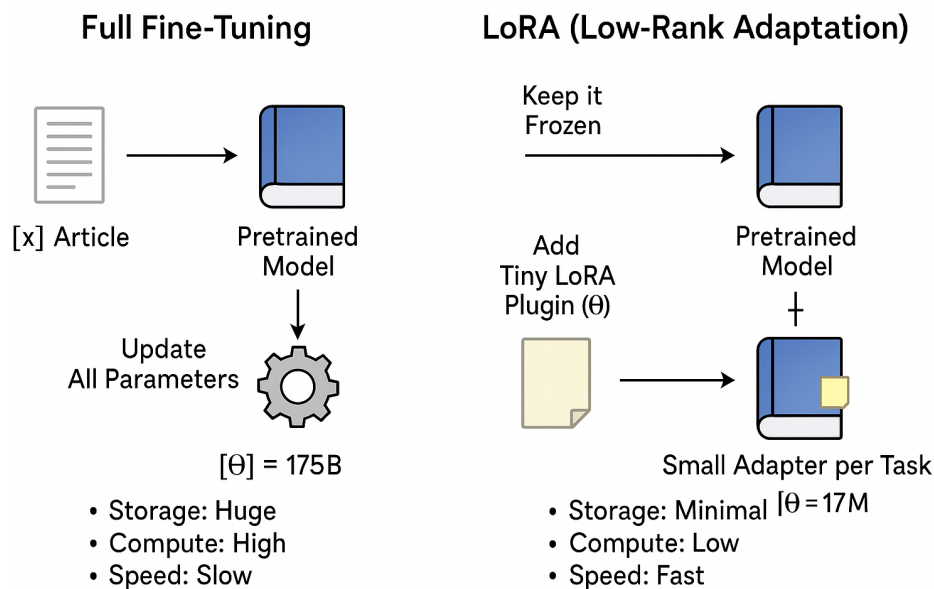


Figure 1: High-level comparison of full fine-tuning and LoRA. Full fine-tuning updates all parameters, while LoRA freezes the base model and trains small adapter modules, greatly reducing compute and storage costs.

## Key Ideas

- Factorize large weight matrices into two smaller low-rank matrices.
- Freeze the pretrained model and update only the new low-rank adapters.
- Maintain high accuracy with significantly fewer trainable parameters.

## Advantages

- Saves GPU memory and storage.
- Enables scalable fine-tuning for multiple downstream tasks.
- Compatible with most transformer-based architectures.

## 2 QLoRA (Quantized LoRA)

While LoRA reduces the number of trainable parameters, very large models (e.g., 65B parameters) still require significant memory to store frozen weights in FP16 precision. QLoRA extends LoRA by combining it with quantization, enabling efficient fine-tuning of massive models on limited hardware.

## Core Ideas

- **4-bit Quantization:** Instead of storing frozen base model weights in FP16, QLoRA compresses them into 4-bit precision, greatly reducing GPU memory usage.
- **LoRA Adapters on Quantized Weights:** LoRA adapters are trained on top of the quantized weights while the original model remains frozen.
- **Double Quantization:** A secondary quantization step ensures that numerical precision errors introduced during 4-bit storage are minimized, preserving accuracy.
- **Paged Optimizers:** Optimizer states are efficiently managed by paging between CPU and GPU memory, enabling training of very large models without exceeding device limits.

## Advantages

- Makes fine-tuning extremely large LLMs (up to 65B parameters) feasible on a single 48GB GPU.
- Maintains accuracy despite aggressive quantization thanks to double quantization and careful optimizer design.
- Reduces VRAM usage to as low as 20GB for fine-tuning multi-billion parameter models.
- Enables democratization of LLM fine-tuning, making advanced research possible on consumer-grade hardware.

## Benchmarks and Results

QLoRA has been shown to achieve performance nearly identical to full-precision fine-tuning across a variety of NLP benchmarks:

- Fine-tuned 65B-parameter models achieve competitive results on benchmarks like MMLU, ARC, and GSM8K.
- Compared to FP16 fine-tuning, QLoRA reduces memory usage by over 75% without significant accuracy loss.
- Enables fine-tuning of 33B models on a single NVIDIA A100 GPU, and 65B models across only a few GPUs, previously considered infeasible.

## Summary

QLoRA pushes the limits of parameter-efficient fine-tuning by merging quantization with LoRA. By storing models in 4-bit precision, applying double quantization, and leveraging paged optimizers, QLoRA enables fine-tuning of massive LLMs on modest hardware. This not only lowers costs but also broadens access to advanced AI research, empowering smaller labs and independent researchers to adapt state-of-the-art models for real-world tasks.

## Conclusion

LoRA and QLoRA represent major breakthroughs in parameter-efficient fine-tuning. **LoRA** introduces low-rank adapters to minimize the number of trainable parameters while freezing the pretrained model. **QLoRA** pushes this idea further by combining LoRA with quantization, enabling fine-tuning of massive LLMs on consumer-grade hardware.

Together, these methods make the adaptation of large models both practical and affordable, accelerating the pace of innovation in natural language processing.