# LoRA and QLoRA: Parameter-Efficient Fine-Tuning of Large Language Models

**Abstract**

This paper presents an overview of two prominent parameter-efficient fine-tuning methods for large language models (LLMs): LoRA (Low-Rank Adaptation) and QLoRA (Quantized LoRA). We discuss their motivations, technical mechanisms, comparative advantages and trade-offs, and practical considerations for deployment.

## 1 Introduction

Fine-tuning very large pretrained language models (LLMs) is computationally demanding both in memory and compute. Traditional full-parameter fine-tuning updates all weights of the model, which becomes infeasible for large models in resource-constrained environments. Parameter-efficient fine-tuning (PEFT) methods reduce the number or size of trainable parameters while preserving performance, enabling practical adaptation of large models. Two influential techniques in this area are LoRA and its more recent extension QLoRA.

## 2 LoRA: Low-Rank Adaptation

### 2.1 Motivation

LoRA was proposed to reduce the cost of fine-tuning large models by only learning a small set of additional adapter parameters, rather than adjusting the full model weights. This makes fine-tuning faster, cheaper, and more flexible.

### 2.2 Advantages & trade-offs

**Advantages**: - Major reduction in number of trainable parameters $\rightarrow$ less GPU memory, faster training, smaller checkpoints. - Modular: different adapters can be maintained for different tasks/applications. - Often close to full-fine-tuning performance for many downstream tasks.

**Trade-offs**: - The expressivity is limited by the low-rank assumption; some tasks may suffer moderate performance drop relative to full fine-tuning. - Choosing appropriate rank $r$ and dropout, learning rate requires care. - While

trainable size is small, inference still uses the full base model (though frozen) so memory/latency for inference remains similar.

# 3 QLoRA: Quantized LoRA

## 3.1 Motivation

As LLMs grow into tens or hundreds of billions of parameters, even freezing the base model but fine-tuning adapters becomes costly in memory (storing the base model) and compute (backprop through huge models). To further reduce resource demands, QLoRA combines LoRA with quantization of the base model weights to very low-bits (e.g., 4-bit) while keeping adapter updates full precision. This enables fine-tuning of massive models on more modest hardware.

## 3.2 Technical mechanism

Key components of QLoRA (from the original work) include: - Quantize the pretrained model weights to 4-bit precision (the paper proposes a data type "NormalFloat4" or NF4). :contentReferenceindex=2 - Freeze the quantized base model weights; insert and train LoRA adapters on top of it. - Use double-quantization and specialized memory optimisers to reduce peak GPU memory usage. :contentReferenceindex=3

By doing so, the authors showed they could fine-tune a 65 B-parameter model on a single 48 GB GPU, achieving 99

## 3.3 Advantages & trade-offs

Advantages: Much lower memory footprint for both fine-tuning and inference (due to quantized base). For example QLoRA uses 75- Enables large scale fine-tuning of very large models on limited hardware. - Good task performance in many cases, often very competitive.

**Trade-offs**: - Quantization may introduce small accuracy degradation or instability. - Training time may be somewhat slower compared to pure LoRA because of quantization/dequantization overhead. :contentReferenceindex=6 - Requires careful engineering (e.g., specialised kernels, memory optimisers) and might be less flexible depending on framework support.

# 4 Comparative Summary

— Method — Trainable ————————————————-———————— —
LoRA — 0.5-5— QLoRA — Similar

A practical guideline: if your base model is manageable and you want fastest training, LoRA is a solid choice. If you have a very large model (e.g., ¿ 10 B parameters) or extreme memory constraints, QLoRA is a powerful approach — provided you are comfortable with quantization trade-offs.

# 5 Practical Considerations

- Frameworks: Both methods are supported in libraries like PEFT (Parameter Efficient Fine Tuning) and Transformers ecosystem. :contentReferenceindex=10 - Adapter checkpointing: With LoRA/QLoRA you typically only save the small adapter weights (A and B matrices) rather than the full model. - Inference/deployment: With QLoRA the base model is quantized, so ensure your runtime supports the quantized model format; testing for accuracy drop is important. - Reproducibility/stability: Some studies show that fine-tuning with QLoRA may exhibit higher variance across runs. :contentReferenceindex=11

# 6 Conclusion

LoRA and QLoRA represent significant advances in adapting very large language models efficiently. LoRA uses a low-rank adaptation strategy to drastically reduce the number of trainable parameters, while QLoRA extends this by combining quantized base models with LoRA adapters to push the limits of model size vs hardware budget. Selecting between them depends on your hardware constraints, model size, and tolerance for engineering complexity. Future work continues to explore even lower-precision quantization, more stable fine-tuning, and domain-specific adaptation for LLMs.