

LoRA and QLoRA - Effective Methods to Fine-tune Your LLMs

Exploring PEFT Approaches for Optimized Large Language Model Fine-Tuning

Author: Mohammad Arshad

Objective: Understand low-resource fine-tuning strategies for high-performing, adaptable NLP models.

Introduction to Fine- tuning LLMs

Definition of Fine-tuning:

- Adjusts parameters in pre-trained language models to adapt to specific NLP tasks.
- Essential for tasks like sentiment analysis, question answering, and language translation.

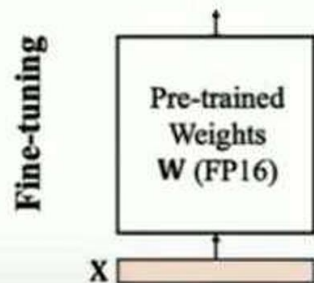
Challenge with Standard Fine-tuning:

- High resource requirements.
- Large model files hinder deployment.

Training Techniques

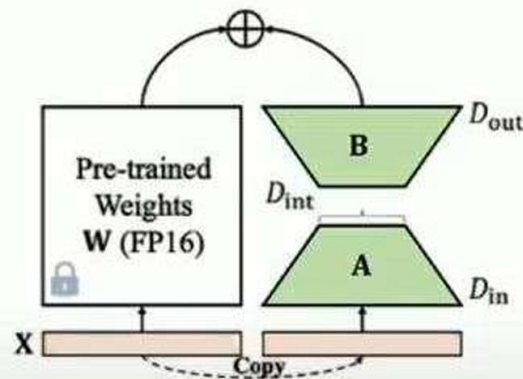
SFT Techniques

Full Fine-Tuning 16-bit precision



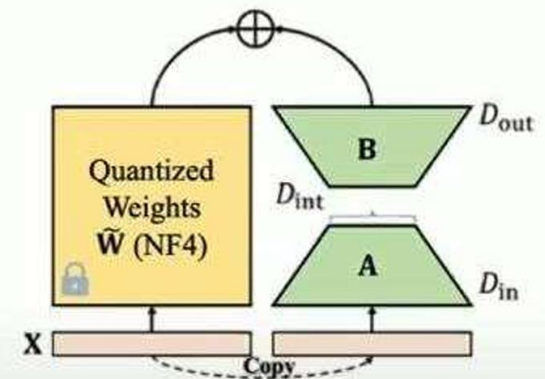
- ✓ Best performance
- ✗ Very high VRAM usage

LoRA 16-bit precision



- ✓ Quick training
- ✗ Still costly

QLoRA 4-bit precision



- ✓ Low VRAM usage
- ✗ Degrades performance

Parameter Efficient Fine- tuning (PEFT)

Overview:

- PEFT minimizes the number of trainable parameters.

Benefits:

- Reduced computation and enhanced portability.

Why Use PEFT?

- Overcomes computational barriers.
- Enhances model portability for deployment in production settings.

PEFT Methods Overview

- LoRA: Optimizes model tuning by adjusting additional parameters.

- QLoRA: Combines LoRA with quantization for efficiency.

Training Techniques

Full Fine-tuning

- Updates all model parameters
- Requires significant computational resources
- Provides maximum adaptation to new tasks
- Can potentially lead to catastrophic forgetting

LoRA (Low-Rank Adaptation)

- Updates a small number of task-specific parameters
- Much more efficient than full fine-tuning
- Preserves most of the original model's knowledge
- Can be combined with the original model weights

QLoRA

- Combines LoRA with quantization techniques (4-bit precision)
- Even more memory-efficient than standard LoRA
- Allows fine-tuning of larger models on consumer hardware
- May have a small trade-off in performance compared to full-precision LoRA



LoRA (Low-Rank Adaptation)

What is LoRA?

- Adds low-rank matrices to fine-tune a subset of parameters.

How It Works:

- 16-bit Transformer, adds extra weights while retaining original ones.

Benefits:

- Retains original model knowledge and reduces trainable parameters.
- 

LoRA Fine- Tuning Process

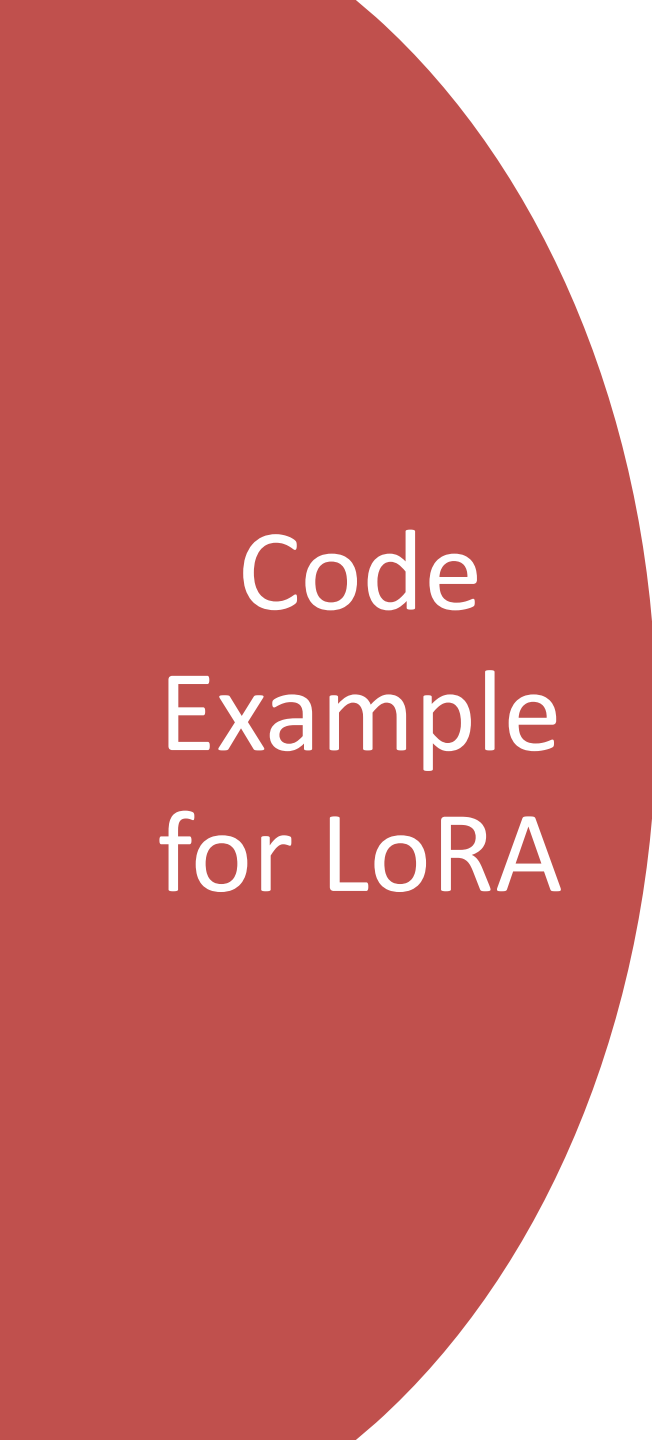
Technical Breakdown:

- Uses a weight matrix (W_0).
- LoRA introduces additional weights.

Advantages:

- Reduces parameters by using low-rank matrices.



A large red circle on the left side of the slide, partially cut off by the edge.


Code Example for LoRA


Objective: Fine-tuning DistilBERT for Sentiment Analysis.

Process:

1. Data Loading
2. Model Setup
3. Tokenization
4. PEFT Configuration
5. Training

Outcome: Reduced resource consumption with competitive performance.

Four purple curved lines of varying lengths and orientations, located in the bottom right corner of the slide.

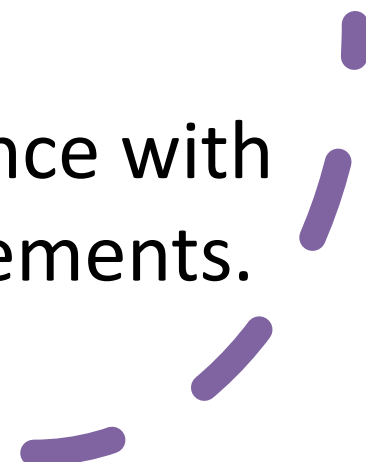


QLoRA (Quantized Low-Rank Adaptation)

Introduction:

- Combines LoRA's low-rank adaptation with quantization.

Why QLoRA?

- Ideal for low-resource devices.
 - Maintains performance with small resource requirements.
- 

Working of QLoRA

Key Components:

- 4-bit Normal Float (NF4):
Efficient storage
- Double Quantization:
Compresses further

Memory Optimization:

- Minimal memory usage with
preserved performance.

QLoRA Process Explained

1. Normalization: Adjusts weights for consistent quantization.
2. Quantization: 4-bit precision storage.
3. Dequantization: Restores weights for accuracy.

Impact: Maintains efficacy with reduced memory.

A large red circle on the left side of the slide, partially cut off by the edge.

Variants of QLoRA

QALoRA:

- Quantizes adapter weights during fine-tuning.

LongLoRA:

- Optimized for extended contexts, ideal for large document tasks.



Summary and Benefits

Summary:

- LoRA: Simplifies fine-tuning.
- QLoRA: Saves memory while retaining performance.

Benefits:

- Efficient resource use, fast adaptation, high portability.

Conclusion and Next Steps

Conclusion:

- LoRA and QLoRA offer scalable, efficient fine-tuning solutions.

Next Steps:

- Experiment with different PEFT techniques for varied tasks.

Q&A: Open floor for questions and discussion.