

LoRA and QLoRA: Parameter-Efficient Fine-Tuning Techniques for LLMs

1 Introduction

The rapid growth of pretrained language models has significantly improved natural language understanding and generation. However, fine-tuning models with billions of parameters remains computationally expensive and memory-intensive, as traditional fine-tuning updates all model weights. To address these challenges, parameter-efficient fine-tuning methods have been proposed. Among the most effective approaches are Low-Rank Adaptation (LoRA) and Quantized Low-Rank Adaptation (QLoRA), which enable efficient adaptation of large models by updating only a small subset of parameters. These methods allow large language models to be fine-tuned on limited hardware while preserving strong performance.

2 Low-Rank Adaptation (LoRA)

LoRA is a fine-tuning technique that reduces the number of trainable parameters by injecting low-rank matrices into pretrained models. Instead of updating the original weight matrices, LoRA learns small rank-decomposed updates that are added to the frozen weights during training.

2.1 Methodology

Given a weight matrix W , LoRA decomposes the update into two low-rank matrices:

$$\Delta W = AB$$

where A and B have much smaller dimensions than W . During training, only A and B are updated, while the original weights remain frozen. This significantly reduces memory usage and computational cost.

2.2 Advantages of LoRA

LoRA offers several benefits:

- **Parameter Efficiency:** Only a small number of parameters are trained.
- **Reduced Memory Usage:** Suitable for limited GPU environments.
- **Scalability:** It can be applied to various transformer-based models like GPT, BERT and T5 making it versatile for different tasks.
- **Modularity:** LoRA adapters can be added or removed without modifying the base model.

2.3 Limitations of LoRA

Despite its efficiency, LoRA still requires storing the base model in full precision, which can be costly for extremely large models.

3 Quantized Low-Rank Adaptation (QLoRA)

QLoRA extends LoRA by introducing model quantization, enabling even larger models to be fine-tuned on consumer-grade hardware. QLoRA quantizes the base model weights to 4-bit precision while keeping LoRA adapters in higher precision. This approach drastically reduces memory consumption without significantly degrading performance.

3.1 Training Strategy

QLoRA fine-tunes large language models efficiently by following a structured process. Each step focuses on reducing memory and computation while adapting the model to a specific task.

1. Quantize the Base Model

- The pretrained model is converted from full precision to 4-bit weights.
- This reduces GPU memory usage, allowing large models to run on smaller hardware.
- Quantization methods such as NF4 help maintain accuracy during compression.

2. Add Low-Rank Adapters

- Small adapter layers are inserted into selected parts of the model, typically the attention layers.
- These adapters remain in higher precision (e.g., 16-bit) to ensure stable training.
- The backbone model is kept frozen, so the original weights are not modified.

3. Fine-Tune Only the Adapters

- During training, only the adapter layers are updated.
- This drastically reduces the number of trainable parameters and the required computation.
- Fine-tuning becomes faster and feasible on a single GPU or low-resource device.

4. Merge or Keep Adapters Separate

- After training, adapters can be merged into the quantized model for deployment.
- Alternatively, they can be kept separate, allowing reuse or swapping for different tasks without retraining the base model.

3.2 Advantages of QLoRA

- **Extreme Memory Efficiency:** Enables training very large models on limited hardware.
- **High Performance:** Achieves performance comparable to full fine-tuning.
- **Scalability:** Makes large-scale experimentation accessible to more users.

4 Applications

LoRA and QLoRA are widely used in:

- Instruction tuning of large language models
- Domain-specific adaptation (medical, legal, financial NLP)
- Conversational agents and chatbots
- Text generation and summarization tasks

5 Conclusion

LoRA and QLoRA represent powerful approaches to efficient fine-tuning of large language models. LoRA provides a lightweight and modular solution, while QLoRA further extends these benefits by enabling low-memory training through quantization. Also, LoRA focuses on reducing trainable parameters, while QLoRA further reduces memory through quantization and enable fine-tuning of larger models compared to LoRA. Together, they significantly lower the barrier to adapting large-scale language models, making advanced NLP systems more accessible.