

LoRTA: LOW RANK TENSOR ADAPTATION OF LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Low Rank Adaptation (LoRA) is a popular Parameter Efficient Fine Tuning (PEFT) method that effectively adapts large pre-trained models for downstream tasks. LoRA parameterizes model updates using low-rank matrices at each layer, significantly reducing the number of trainable parameters and, consequently, resource requirements during fine-tuning. However, the lower bound on the number of trainable parameters remains high due to the use of the low-rank matrix model. In this paper, we address this limitation by proposing a novel approach that employs a low rank tensor parameterization for model updates. The proposed low rank tensor model can significantly reduce the number of trainable parameters, while also allowing for finer-grained control over adapter size. Our experiments on Natural Language Understanding, Instruction Tuning, Preference Optimization and Protein Folding benchmarks demonstrate that our method is both efficient and effective for fine-tuning large language models, achieving a reduction in the number of parameters while maintaining comparable performance.

1 INTRODUCTION

The advent of Large Language Models (LLMs) – billion parameter scale models pre-trained on vast corpora of data – has enabled unprecedented capabilities across a wide range of tasks. However, as LLM sizes continue to grow exponentially, their computational and memory demands represent significant challenges, particularly for those lacking access to high-performance computing infrastructure (Varoquaux et al., 2024). This has spurred interest in parameter efficient fine tuning (PEFT) techniques (Ding et al., 2023), which facilitate the adaptation of LLMs to specific applications, downstream tasks or user preferences, by using only a small fraction of trainable parameters. Most importantly, they reduce GPU memory requirements, primarily by shrinking optimizer states (Liao et al., 2023). Moreover, they provide greater efficiency in storage and deployment, enabling the management of multiple fine-tuned LLMs with reduced storage footprints and faster load times (Sheng et al., 2023; Wen & Chaudhuri, 2024), which is particularly relevant for applications requiring rapid model switching across numerous task- or user-specific models.

Beyond computational benefits, PEFT techniques can also mitigate overfitting risks associated with fine-tuning high-capacity LLMs. By constraining model updates, PEFT methods can act as an implicit regularization mechanism, improving generalization (Fu et al., 2023; Sun et al., 2023). Parameter sharing, a well-established technique in deep learning architecture design, has been shown to improve generalization across various tasks such as protein folding (Jumper et al., 2021; Lin et al., 2023), image segmentation (Ronneberger et al., 2015), and generative modeling (Rombach et al., 2022). Incorporating parameter sharing in PEFT methods has also improved performance in specialized applications with limited data, such as in medical domains (Dutt et al., 2023; Zhu et al., 2024).

Low Rank Adaptation (LoRA) is a popular PEFT approach that uses a low rank parameterization of weight matrix updates (Hu et al., 2021). For instance, these allow to fine tune a 175 billion parameter LLM using only 5 million trainable parameters (Hu et al., 2021) without performance degradation. Since the model updates can be merged with the frozen weights, LoRA incurs no additional inference cost when deployed, unlike prompt (Li & Liang, 2021a; Liu et al., 2023) and adapter-based (Houlsby et al., 2019; He et al., 2021; Pfeiffer et al., 2020) PEFT methods. However, the lower bound on trainable parameters often remains substantial for large-scale models. Recent works have aimed to further reduce the number of parameters in LoRA by allocating different ranks

across update matrices in different layers (Valipour et al., 2022; Zhang et al., 2023b), using fixed low rank projections (Zhang et al., 2023a; Kopiczko et al., 2023), and parameterizing low rank matrices using a random basis (Koohpayegani et al., 2024).

In this work, we introduce Low Rank Tensor Adapters (LoRTA), which exploit redundancy in weight updates across different layers, heads, and attention matrices by representing updates as a unified 5th-order low-rank tensor model. This holistic approach not only reduces the number of trainable parameters but also facilitates learning by exploiting shared information among various model components. We show that the parameter-sharing schemes in state-of-the-art methods (Kopiczko et al., 2023; Koohpayegani et al., 2024) are implicit low-rank tensor models with fixed random factors, while our explicit tensor model is *fully trainable* and more expressive. Our higher-order tensor update and Candecomp-Parafac model (Harshman & Lundy, 1994) enjoys greater parameter efficiency and favorable scaling over low-rank tensor models recently proposed for vision transformers (Jie & Deng, 2023; Edalati et al., 2023) and LLMs (Yang et al., 2024; Bershtsky et al., 2024).

The two main advantages of our proposed method over existing matrix and tensor based approaches can be summarized as follows:

- (A1) It enables a reduction in the number of trainable parameters by using updates with lower tensor rank.
- (A2) The number of parameters scales better with tensor rank, number of fine-tuned matrices, embedding dimension and attention heads.

As a consequence, our approach also provides finer-grained control of adapter size. We evaluate our method on diverse benchmarks including Natural Language Understanding, Instruction Tuning, Preference Optimization, and Protein Folding. Our experiments demonstrate that LoRTA can achieve up to an order of magnitude reduction in the number of trainable parameters compared to state-of-the-art PEFT methods, with minimal performance trade-offs.

2 PRELIMINARIES

2.1 TRANSFORMER ARCHITECTURE

We focus on the transformer architecture, although it can be naturally extended to other architectures such as Convolutional Neural Networks and Long Short Term Memory networks. We adopt the problem setting presented in (Thickstun, 2021). In the transformer model, an initial embedding layer maps input tokens to d -dimensional vector representations. These embeddings then pass through a series of layers, each performing multi-head attention, normalization and feed-forward operations. The input to the l -th layer of the transformer is a matrix $\mathbf{X}^{(l)} \in \mathbb{R}^{N \times d}$, where N is the number of queries, represented in a d -dimensional feature space. A vanilla transformer layer with H attention heads is then defined as follows:

$$\mathbf{X}^{(l+1)} = \text{LayerNorm} \left(\mathbf{Y}^{(l)} + \text{MLP} \left(\mathbf{Y}^{(l)} \right) \right) \quad (1)$$

$$\mathbf{Y}^{(l)} = \text{LayerNorm} \left(\mathbf{X}^{(l)} + \text{Attn} \left(\mathbf{X}^{(l)} \right) \right) \quad (2)$$

$$\text{Attn} \left(\mathbf{X}^{(l)} \right) = \mathbf{X}^{(l)} + \sum_{h=1}^H \text{softmax} \left(\frac{\mathbf{X}^{(l)} \mathbf{Q}_h^{(l)} \mathbf{K}_h^{(l)T} \mathbf{X}^{(l)T}}{\sqrt{d}} \right) \mathbf{X}^{(l)} \mathbf{V}_h^{(l)} \mathbf{P}_h^{(l)T} \quad (3)$$

$$\text{MLP} \left(\mathbf{X}^{(l)} \right) = \text{ReLU} \left(\mathbf{X}^{(l)} \mathbf{G}_1^T + \mathbf{1}_N \mathbf{b}_1^T \right) \mathbf{G}_2^T + \mathbf{1}_N \mathbf{b}_2^T, \quad (4)$$

where $\mathbf{K}_h^{(l)}, \mathbf{Q}_h^{(l)}, \mathbf{V}_h^{(l)}, \mathbf{P}_h^{(l)} \in \mathbb{R}^{d \times d_H}$ are the key, query, value and projection matrices respectively, for head h and layer l .

2.2 LOW RANK (MATRIX) ADAPTATION

LoRA modifies the pre-trained weights by adding a trainable update. Explicitly, at each layer and head h :

$$\mathbf{K}_h = \mathbf{K}_h^0 + d\mathbf{K}_h, \quad \mathbf{Q}_h = \mathbf{Q}_h^0 + d\mathbf{Q}_h, \quad \mathbf{V}_h = \mathbf{V}_h^0 + d\mathbf{V}_h, \quad \mathbf{P}_h = \mathbf{P}_h^0 + d\mathbf{P}_h, \quad (5)$$

where $\mathbf{K}^0, \mathbf{Q}^0, \mathbf{V}^0, \mathbf{P}^0$ denote the pre-trained weights and $d\mathbf{K}, d\mathbf{Q}, d\mathbf{V}, d\mathbf{P}$ the trainable adapters.

While each attention head’s MLP contains two trainable matrices, \mathbf{G}_1 and \mathbf{G}_2 , our focus is on fine-tuning the attention matrices. This has been demonstrated to be effective for LLM adaptation (Hu et al., 2021; Kopiczko et al., 2023). Nevertheless, these methods can be easily extended to other parameters, including the MLP weights.

Let $\mathbf{W}_h \in \{\mathbf{Q}_h, \mathbf{K}_h, \mathbf{V}_h, \mathbf{P}_h\}$ for $h = 1, \dots, H$ denote the query, key, value and projection matrices, respectively, for each attention head. After concatenating updates across all attention heads, we get:

$$d\tilde{\mathbf{W}} = (d\mathbf{W}_1, \dots, d\mathbf{W}_H) \in \mathbb{R}^{d \times d}.$$

Hu et al. (2021) proposed to parametrize the updates using rank- r matrices, which can be expressed as

$$d\tilde{\mathbf{W}} = \frac{\alpha}{r} \mathbf{A} \mathbf{B}^T, \quad \mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times r}, \quad (6)$$

where α is a constant and r denotes the rank of the update. The scaling factor simply aims to reduce the efforts of re-tuning the learning rate when training adapters of varying rank. It has been shown that while this scaling heuristic works well for smaller ranks, it can be sub-optimal for larger ranks (Kalajdziewski, 2023). Hayou et al. (2024) have also shown that setting the learning rate for the \mathbf{A} and \mathbf{B} matrices appropriately can further improve convergence and performance.

Although LoRA is an efficient fine-tuning technique, the number of parameters required for each layer is at least $8 \cdot d \cdot r$. Thus, the total number of trainable parameters is:

$$\# \text{parameters (LoRA)} = 8 \cdot d \cdot L \cdot r, \quad (7)$$

where L is the total number of layers. Even with $r = 1$, this results in $8 \cdot d \cdot L$ parameters. In practice, for LLMs with high dimensionality (d) and many layers (L), this lower bound can still lead to a significant number of trainable parameters.

LoRA has also been combined with model weight quantization (Dettmers et al., 2024), further decreasing resource requirements. Unlike adapter-based PEFT methods (Houlsby et al., 2019; Pfeiffer et al., 2020; Rücklé et al., 2020; He et al., 2021), LoRA does not introduce additional inference time overhead during deployment, as the trainable matrices can be integrated with the fixed weights.

Building upon this foundation, AdaLoRA (Zhang et al., 2023b) expands the LoRA technique by introducing dynamic rank adjustment for low-rank matrices during fine-tuning. The fundamental concept involves optimally allocating the parameter resources by selectively pruning less crucial components of the matrices based on an importance metric. LoRA-FA (Zhang et al., 2023a) reduces the number of trainable parameters by freezing the \mathbf{A} matrix to its random initialisation, while achieving similar performance to LoRA.

2.3 TENSOR ALGEBRA

In the following sections we introduce our proposed LoRTA framework, which is a tensor adaptation model for PEFT. To facilitate the upcoming analysis, we briefly present some tensor algebra preliminaries and refer the reader to Appendix A and Sidiropoulos et al. (2017); Kolda & Bader (2009) for further details.

A N -order tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ is an N -way array indexed by i_1, i_2, \dots, i_N with elements $\mathcal{X}(i_1, i_2, \dots, i_N)$. It consists of N types of modes: $\mathcal{X}(:, i_2, \dots, i_N)$, $\mathcal{X}(i_1, :, \dots, i_N)$, \dots , $\mathcal{X}(i_1, i_2, \dots, :)$. Any tensor can be decomposed as a sum of N -way outer products as:

$$\mathcal{X} = \sum_{r=1}^R \mathbf{A}_1[:, r] \circ \mathbf{A}_2[:, r] \circ \dots \circ \mathbf{A}_N[:, r], \quad (8)$$

where $\mathbf{A}_n = [\mathbf{a}_n^1, \mathbf{a}_n^2, \dots, \mathbf{a}_n^R] \in \mathbb{R}^{I_n \times R}$, $n = 1, \dots, N$ are called the low rank factors of the tensor. The above expression represents the *canonical polyadic decomposition* (CPD) or *parallel factor analysis* (PARAFAC) (Harshman & Lundy, 1994) of a tensor. A tensor can be fully characterized by its latent factors, so we can represent a tensor by its CPD model as:

$$\mathcal{X} = \llbracket \mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_N \rrbracket. \quad (9)$$

Unlike other tensor models, such as Tucker and Block Term Decomposition (BTD), the CPD model is unique under certain conditions. As a result, the CPD model is often preferred when the goal is to minimize the number of parameters.

A tensor can also be represented as a set of matrices, by fixing all the modes but two as:

$$\mathcal{X}[:, :, i_3, \dots, i_N] = \mathbf{A}_1 (\text{Diag}(\mathbf{A}_3(i_3, :)) \odot \dots \odot \text{Diag}(\mathbf{A}_N(i_N, :))) \mathbf{A}_2^T, \quad (10)$$

where $\text{Diag}(\mathbf{A}_n(i_n, :))$ is the diagonal matrix with diagonal equal to $\mathbf{A}_n(i_n, :)$.

3 LOW RANK TENSOR ADAPTATION

3.1 PARAMETER SHARING ACROSS LAYERS

To further increase the compression ratio in PEFT models, recent works (Kopiczko et al., 2023; Koohpayegani et al., 2024) suggest sharing parameters across layers that operate as predefined projection matrices. As we see next, this leads to tensor factorization models with fixed parameters.

Vector-based Random Matrix Adaptation (VeRA) Kopiczko et al. (2023) have proposed to parameterize updates using two learnable vectors at each layer and fixed random matrices shared across all layers. The update at each layer can be expressed as

$$d\tilde{\mathbf{W}} = \mathbf{A} \text{Diag}(\mathbf{C}_D[l, :]) \mathbf{B}^T \text{Diag}(\mathbf{C}_B[l, :]), \quad (11)$$

where $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times r}$ are the random projections, and $\mathbf{C}_D \in \mathbb{R}^{L \times r}$, $\mathbf{C}_B \in \mathbb{R}^{L \times d}$ are matrices that collect trainable vectors across layers. The model in 11 is a coupled matrix factorization model and is similar to a tensor model. In particular, if we remove the \mathbf{C}_B term VeRA can be interpreted as a low-rank tensor CPD parameterization with fixed random factors. That is, the weight update $\tilde{\mathbf{W}}$ is a rank- r third order tensor $\mathcal{T} \in \mathbb{R}^{d \times d \times L}$. Note that, omitting the \mathbf{C}_B term has been shown to lead to a small performance degradation unlike omitting \mathbf{C}_D (Kopiczko et al., 2023).

Random Matrix basis Adaptation (NOLA) In a similar manner, Koohpayegani et al. (2024) have proposed to parameterize the weight update by expressing the matrices \mathbf{A} and \mathbf{B} as linear combinations of fixed random basis matrices, that are shared across all layers. The weight update $d\tilde{\mathbf{W}}$ for layer l is then given by:

$$d\tilde{\mathbf{W}}_l = \sum_{i=1}^k \sum_{j=1}^k \alpha_{(i,l)} \beta_{(j,l)} \mathbf{A}_i \mathbf{B}_j^T, \quad (12)$$

where $\mathbf{A}_i, \mathbf{B}_j \in \mathbb{R}^{d \times r}$ are fixed random matrices, shared across all layers, and $\alpha_l = \{\alpha_{(i,l)}\}_{i=1}^K$ and $\beta_l = \{\beta_{(j,l)}\}_{j=1}^K$ are the learned coefficients for each layer. If we stack the random matrices $\mathbf{A}_i, \mathbf{B}_j \in \mathbb{R}^{d \times r}$ into tensors \mathcal{A}, \mathcal{B} such that: $\mathcal{A}[:, :, i] = \mathbf{A}_i$ and $\mathcal{B}[:, :, j] = \mathbf{B}_j$, then 12 can be cast as:

$$d\tilde{\mathbf{W}}_l = \sum_{i=1}^k \sum_{j=1}^k \alpha_{(i,l)} \beta_{(j,l)} \sum_{m=1}^r \mathcal{A}[:, m, i] \mathcal{B}[:, m, j]^T = \sum_{m=1}^r \mathcal{A}[:, m, :] (\alpha_l \beta_l^T) \mathcal{B}[:, m, :]^T, \quad (13)$$

and $d\tilde{\mathbf{W}}_l$ admits the following factorization. $d\tilde{\mathbf{W}}_l = \sum_{m=1}^r \mathbf{P}_A^{(m)} (\alpha_l \beta_l^T) \mathbf{P}_B^{(m)T}$, where $\mathbf{P}_A^{(m)} = \mathcal{A}[:, m, :]$, and $\mathbf{P}_B^{(m)} = \mathcal{B}[:, m, :]$ are also random projection matrices with different dimensions compared to $\mathbf{A}_i, \mathbf{B}_j$. As a result, NOLA can be viewed as the following tensor factorization model:

$$d\tilde{\mathcal{W}} = \sum_{m=1}^r \left[\mathbf{P}_A^{(m)} \tilde{\mathbf{A}}, \mathbf{P}_B^{(m)} \tilde{\mathbf{B}}, \mathbf{I} \right], \tilde{\mathbf{A}}[:, l] = \alpha_l, \tilde{\mathbf{B}}[:, l] = \beta_l. \quad (14)$$

The expression in 14 is a summation of CPD models, also known as Block Term Decomposition, which is an expressive tensor model, but can lack parsimony (Kolda & Bader, 2009).

3.2 LORTA: A MORE EFFICIENT TENSOR MODEL

In the previous section, we explored PEFT models that share parameters across layers, highlighting their correspondence to tensor factorization models. Namely, VeRA and NOLA utilize fixed projection matrices shared across layers. However, this strategy can result in models that are larger than necessary relative to their degrees of freedom due to the inclusion of these random matrices. Although these matrices can be generated on the fly by solely storing the pseudo-random number generator seed, this still incurs additional resource demands during training, and increases loading time for inference.

To address this issue, we propose modeling the trainable adapters using a low-rank CPD structure. This choice is motivated by the favorable properties of CPD: it is universal, capable of exactly factorizing any tensor, yet remains concise and parsimonious, typically requiring only a small number of parameters to achieve low approximation error (Sidiropoulos et al., 2017). This contrasts with tensor adapters used in vision (Jie & Deng, 2023) and recently in LLM finetuning (Bershtsky et al., 2024), which rely on Tucker and Tensor-Train models. In fact, for small ranks, CPD is equivalent to Tucker when the core tensor in Tucker is the identity tensor. However Tucker is always parametrized with a dense tensor and therefore requires a larger number of parameters for the same rank.

LoRTA represents all weight updates as a 5th-order tensor $d\tilde{\mathcal{W}} \in \mathbb{R}^{d \times \frac{d}{H} \times H \times L \times M}$. By integrating updates of layers, heads and the $\mathbf{Q}, \mathbf{K}, \mathbf{V}, \mathbf{P}$ matrices into a unified low-rank CPD tensor model, LoRTA exploits redundancy across different modes of the tensor. This approach can thus not only improve parameter efficiency but also facilitate learning by exploiting the shared information among various components of the model. This contrasts with existing PEFT approaches, which tensorize each weight update independently (Yang et al., 2024) or only share parameters across layers (Jie & Deng, 2023; Bershtsky et al., 2024). In order to illustrate how additional tensor modes can result in parameter efficiency gains, Figure 1 compares – for a single weight update – LoRA with a rank one tensor model that adds attention heads as a mode.

By utilizing a low-rank CPD model, we can express this tensor as:

$$d\tilde{\mathcal{W}} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C}_H, \mathbf{C}_L, \mathbf{C}_M \rrbracket,$$

where $\mathbf{A} \in \mathbb{R}^{d \times r}$ and $\mathbf{B} \in \mathbb{R}^{\frac{d}{H} \times r}$ are factor matrices for the input and output dimensions, respectively, and $\mathbf{C}_H \in \mathbb{R}^{H \times r}$, $\mathbf{C}_L \in \mathbb{R}^{L \times r}$, $\mathbf{C}_M \in \mathbb{R}^{4 \times r}$ are factor matrices for the attention heads, layers, and the four matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V}, \mathbf{P}$. Each weight matrix update can then be retrieved as:

$$d\tilde{\mathcal{W}}[:, :, k, l, i] = \mathbf{A} (\text{Diag}(\mathbf{C}_H[k, :]) \text{Diag}(\mathbf{C}_L[l, :]) \text{Diag}(\mathbf{C}_M[i, :])) \mathbf{B}^\top,$$

where k indexes the attention heads, l indexes the layers, and i indexes the matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V}, \mathbf{P}$. Note that, unlike previous implicit tensor models such as NOLA and VeRA, which rely on fixed random projections or parameters and learn only scaling coefficients, our proposed model is *fully trainable*. All factor matrices ($\mathbf{A}, \mathbf{B}, \mathbf{C}_H, \mathbf{C}_L, \mathbf{C}_M$) are learned during training, providing greater expressiveness and forgoing the dependency on pre-defined random bases or projections.

Table 1 shows how the CP low rank tensor parameterization leads to better scaling in the number of parameters with respect to the tensor rank r . Moreover, our higher-order weight update tensorization improves scaling in terms of transformer architecture hyperparameters, namely the embedding dimension d , number of attention heads H , and number of fine-tuned attention matrices M .

3.3 OTHER LOW RANK TENSOR MODELS IN PEFT

As mentioned in the previous section, existing PEFT tensor-based models differ from our method both in their parameter-sharing schemes, which result from different weight update tensorization approaches, as well as in the low-rank tensor models they employ. Below, we provide a concise overview of these approaches which intends to highlight the differences with LoRTA; further details are available in Appendix X and the provided references.

FaCT & LoTR In the context of vision transformers, Jie & Deng (2023) have proposed to represent updates across all layers as a third order tensor $d\tilde{\mathcal{W}} \in \mathbb{R}^{L \times d \times d}$. They propose two parameterizations of $d\tilde{\mathcal{W}}$, namely, a Tensor Train and Tucker3 low rank tensor models. Recently, Bershtsky et al.

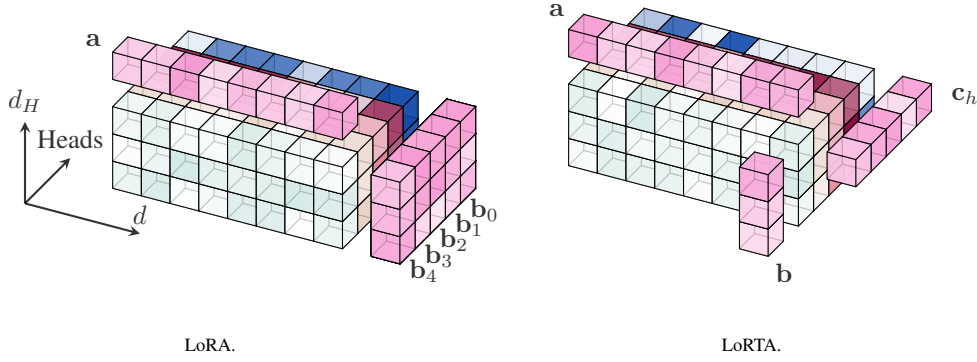


Figure 1. Illustration of a rank 1 adapter for a single weight matrix with multiple heads. (Left) The LoRA update for head h is computed as $d\mathbf{W}_h = \mathbf{b}_h \circ \mathbf{a}$. (Right) The update using a third order low rank tensor model is computed as $d\mathbf{W}_h = \mathbf{b} \circ \mathbf{c}[h] \circ \mathbf{a}$. Both models have the same tensor rank, but the latter has less parameters.

(2024) have proposed to apply the same tensorization across layers to fine-tune LLMs, but using a low rank Tucker2 tensor model to parameterize updates.

LoreTTA Yang et al. (2024) propose two methods that employ low rank tensor models. However, these models do not share parameters across layers, they reparameterize low rank matrix adapters using low rank tensor models. In **LoreTTA-rep** a low rank matrix model is first applied to each weight update in the same manner as described for LoRA in Equation (6). Then each of the ML resulting LoRA factors $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times r}$ are expressed as a n -th order tensor with arbitrary dimensions, i.e. $\mathcal{A}, \mathcal{B} \in \mathbb{R}^{k_1 \times \dots \times k_N}$. Finally, each of these tensors is parametrized Tensor Train model, explicitly, $\mathcal{A} = \prod_{i=1}^N \mathbf{G}_i$ where $\mathbf{G}_i \in \mathbb{R}^{r \times k_i \times r}$. We highlight that the added dimensions k_i are hyperparameters that must satisfy $\prod_i k_i = dr$ and $k_i \geq r$ for all i ; otherwise, it would induce a new tensor rank deficiency. Moreover, choosing appropriate values of k_i might be challenging and necessitate further hyperparameter tuning. Yang et al. (2024) also proposed **LoReTTA-adp**, applying a tucker parameterization to an adapter method, which unlike our method and the rest of the aforementioned methods adds new parameters to the model and thus can not be merged into the original weights, thereby incurring additional inference costs.

4 EXPERIMENTS

4.1 NATURAL LANGUAGE UNDERSTANDING

We evaluate our approach by fine-tuning RoBERTa models on the General Language Understanding Evaluation (GLUE) Wang et al. (2018) benchmark. We conduct experiments across three distinct settings previously reported in the literature by Bershtsky et al. (2024), Yang et al. (2024) and Kopiczko et al. (2023). These settings differ in hyperparameters, including the number of training

Method	Update Tensor shape	Tensor Model	Parameters	r=4	r=64
LoRA	$ML \times d \times d$	Matrix-Batch	$2MLdr$	2.1M	33M
LoReTTA	$ML \times k_1 \times \dots \times k_6$	Custom	$MLr^2 \sum_i k_i$	92k	50M
LoTR	$ML \times d \times d$	Tucker2	$MLr^2 + 2dr$	33k	786k
FacT-TT	$ML \times d \times d$	Tensor-Train	$MLr^2 + 2dr$	33k	786k
FacT-TK	$ML \times d \times d$	Tucker3	$(2d + ML)r + r^3$	33k	790k
Ours	$M \times L \times d \times d/h \times h$	CP	$(d + d/h + h + L + M)r$	17k	274k

Table 1: Number of parameters of different Tensor based PEFT methods as a function of the number of finetuned attention/projection matrices M , the number of layers, L , the embedding dimension d , the number of heads h and the tensor rank of the update, r . For LoreTTA, k_i are hyperparameters that must satisfy $\prod_i k_i = dr$ and $k_i \geq r$ for all i . We also include the number of parameters for the Llama2-7b architecture when finetuning only $M=2$ attention matrices (e.g. Q and V) for different ranks. For LoreTTa we use $k_1 = \dots = k_6 = 5$ for $r = 4$ and $k_1 = k_2 = k_3 = 64$ for $r = 64$.

epochs, different learning rates for the classification head and encoder, the learning rate decay strategy (linear vs fixed), the use of different scaling parameters α , and the grid search ranges. Because best results on the validation set are reported, performance for the same baseline method can vary considerably across settings (see, for example, LoRA performance reported by Hu et al. (2021), Yang et al. (2024) and Bershatsky et al. (2024)). Therefore, we provide an evaluation of our method in a variety of experimental conditions, while also maintaining the original configurations in which state-of-the-art methods were originally reported. Detailed settings can be found in Table ?? in Appendix E.1.

We also finetuned Llama2 models (Touvron et al., 2023) on question-answering (QA) tasks SQuAD (Rajpurkar et al., 2016), DROP (Dua et al., 2019), COPA (Roemmele et al., 2011), and ReCoRD (Zhang et al., 2018), following the experimental setting outlined by Yang et al. (2024). For these tasks, we used a randomly selected subset of 1,000 samples to simulate a low-data regime and increase the task difficulty. All classification tasks are tackled as language modeling tasks following the prompt-based fine-tuning approach described by Malladi et al. (2023).

Baselines We benchmark our method against the following methods:

- **Full finetuning**: all parameters are trained.
- **IA3** (Liu et al., 2022): rescales activations with learned vectors
- **Prefix** (Li & Liang, 2021b): prepends learnable continuous vectors (prefixes) to the input embeddings.
- **LoRA** (Hu et al., 2021), **LoRA-FA** (Zhang et al., 2023a) and **VeRA** (Kopiczko et al., 2023), **LoTR** (Bershatsky et al., 2024), **LoReTTA** (Yang et al., 2024): As previously described.
- We omit **Adapter^H** (Houlsby et al., 2019), **Adapter^P** (Pfeiffer et al., 2020), **Bitfit** (Zaken et al., 2021), **AdapterDrop** (Rücklé et al., 2020), and other methods that are customarily reported but have been outperformed by more recent methods in these settings.

The results in Table 2 show that LoRTA can achieve comparable or slightly superior performance with less trainable parameters compared to state of the art tensor based PEFT methods LoreTTA (Yang et al., 2024) and LoTR (Bershatsky et al., 2024) when finetuning RoBERTa base on GLUE tasks. Similarly, for RoBERTa large LoRTA can also achieve a 6x reduction in the number of trainable parameters with only small drop in average performance (2%) when compared to Kopiczko et al. (2023). In this settings we did not tune the hyperparameters for our method as extensively as baselines, and thus this gap could be further reduced.

In Llama QA experiments, shown in Table 3, full fine-tuning (FT) achieves the highest average score (77.3) with 7 billion trainable parameters, but among the PEFT methods LoRTA (r=8) achieves the highest average score (76.7) with just 0.03 million parameters, representing a 17x reduction in parameter count with respect to the most efficient method.

4.2 INSTRUCTION TUNING

We fine-tune the 7 billion parameter Llama2 (Touvron et al., 2023) models on the cleaned Alpaca instruction tuning dataset (Taori et al., 2023). We train for one epoch, preceded by a warm-up learning rate sweep as in the standard setting. Other hyperparameters are detailed in Appendix E.2.

As shown in Figure 2, LoRTA effectively reduces the number of parameters to a fraction of those required by the lowest rank in LoRA, with only a small performance cost. The loss decreases monotonically with the number of parameters used, both in training and testing, and LoRTA even demonstrates superior performance with fewer parameters for ranks 96 and 192. To further evaluate the fine-tuned models, we use MT-Bench (Zheng et al., 2023), an LLM-as-a-judge benchmark. MT-Bench assesses multi-turn conversational and instruction-following abilities on 80 open-ended questions, covering diverse capabilities such as roleplaying, reasoning, coding and information retrieval. GPT-4 is used to score the outputs of the model on a scale of one to ten.

As shown in Figure 3, LoRTA can almost match average performance despite using just 1/5th of the parameters (r=48). Unlike the loss observed in the Alpaca dataset, performance does not increase monotonically, potentially due to overfitting. Moreover, performance varies across tasks. For example, most LoRTA models surpass LoRA in reasoning but fall short in writing.

	Method	# Trainable Parameters	SST-2	MRPC	CoLA	QNLI	RTE	STS-B	Avg.
LoReTTA	LoRA (r=8)	630k	94.01	91.48	62.08	92.39	74.51	84.69	83.19
	LoReTTA rep	70k	94.28	90.63	61.72	92.40	74.42	89.24	83.78
	LoRTA (r=20)	48k	94.27	92.04	63.35	91.48	75.09	89.82	84.34
	LoRTA (r=12)	29k	93.81	91.13	61.40	92.04	74.73	89.64	83.79
LoTR	LoRA (r=8)	300k	94.2	88.0	61.1	91.3	73.0	90.7	83.05
	LoTR	74k	93.0	85.9	60.5	90.0	66.0	91.9	81.22
	LoRTA (r=16)	15k	94.73	90.44	64.32	92.37	76.9	90.25	84.84
	LoRTA (r=4)	3.4k	94.61	89.21	60.55	90.61	76.9	89.97	83.6
VeRA	LoRA	800k	96.2	90.2	68.2	94.8	85.2	92.3	87.8
	LoRA-FA	3.7M	96.0	90.0	68.0	94.4	86.1	92.0	87.7
	VeRA	61k	96.1	90.9	68.0	94.4	85.9	91.7	87.8
	LoRTA (r=8)	9k	96.3	89.5	65.1	94.3	85.6	91.1	85.7

Table 2: Performance of RoBERTa Base and Large models on GLUE tasks under three different experimental settings reported by LoReTTA (Yang et al., 2024), LoTR (Bershtatsky et al., 2024), and VeRA (Kopiczko et al., 2023). In LoReTTA, LoRTA is applied to the encoder and LoRA to the classifier with the same rank, while for LoTR and VeRA, LoRTA is applied only to the encoder. Trainable parameters include the classifier for LoReTTA but exclude it for LoTR and VeRA, where it is fully trained. VeRA results use RoBERTa Large, whereas LoTR and LoReTTA use RoBERTa Base

Method	# Trainable Parameters	COPA	ReCoRD	SQuAD	DROP	Avg.
FT	7B	86	81.1	90.71	51.38	77.3
LoRA (r=8)	4.19M	81	79.4	90.56	45.96	74.2
Prefix	1.31M	83	81.0	90.56	45.95	75.1
IA3	0.60M	80	81.5	89.41	39.37	72.6
LoReTTA rep	0.51M	86	80.3	88.47	42.71	74.4
LoRTA (r=4)	0.02M	87	81.1	87.4	44.04	74.9
LoRTA (r=8)	0.03M	87	81.6	88.5	49.7	76.7

Table 3: LLama2-7B performance on SuperGLUE and question-answering tasks (SQuAD, DROP). We follow the experimental setup used by Yang et al. (2024).

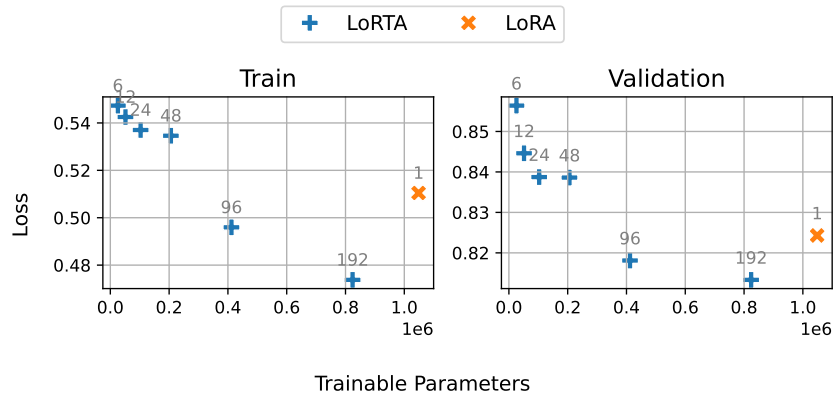


Figure 2. Mean cross-entropy loss on training and testing data for Llama2-7b on the Alpaca dataset vs number of trainable parameters for different adapter ranks. Lower is better. Numbers on top of markers denote the adapter rank.

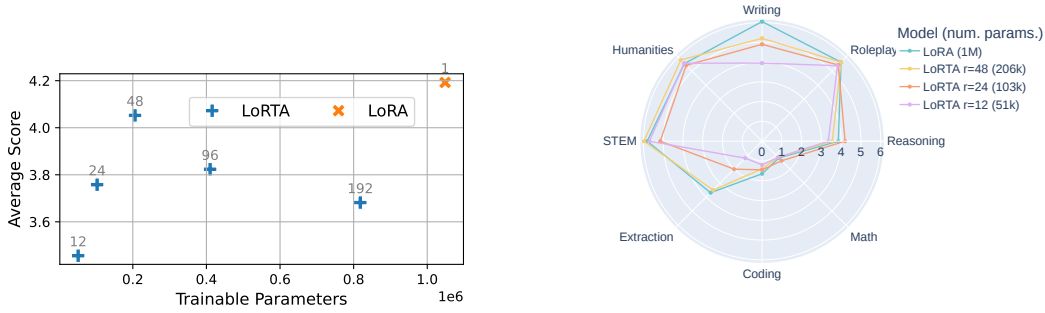


Figure 3. Performance on MT-Bench (Zheng et al., 2023) for Llama2-7b (Touvron et al., 2023) models fine-tuned with LoRA and LoRTA. Higher is better. **Left:** Average score across all questions vs number of trainable parameters. Numbers on top of markers denote the adapter rank. **Right:** Average score per task.

4.3 PREFERENCE OPTIMIZATION

While various techniques to align LLMs with human preferences on specific tasks exist (see, for example, Kaufmann et al. (2023) and references therein), we utilize Direct Preference Optimization (DPO) (Rafailov et al., 2024) due to its simplicity and effectiveness.

We fine-tune the 7 billion parameter Llama 2 model (Touvron et al., 2023) on the cleaned version of the Intel Orca dpo pairs dataset¹. This synthetic preference dataset comprises 6k prompts across various domains and tasks, along with the corresponding outputs from ChatGPT and Llama2-13B. In this version of the dataset, ChatGPT is used to score outputs and the preferred choices are designated based on these scores. Because preference datasets are often small, a KL regularization that penalizes deviations from the pre-trained model’s outputs is used to mitigate overfitting. In our experiments, the regularization coefficient β was set to 0.1. We use Huggingface Transformer Reinforcement Learning (trl) library². For a complete description of hyperparameters see Appendix E.3.

As shown in Figure 4 LoRTA exhibited non-monotonic performance across ranks. This suggests that further hyperparameter tuning may be necessary to stabilize its performance. Although we did not tune hyperparameters, most ranks still outperformed LoRA with significantly fewer parameters. We further evaluated the fine-tuned models on the LLM-as-a-judge MT-benchmark. In this setting,

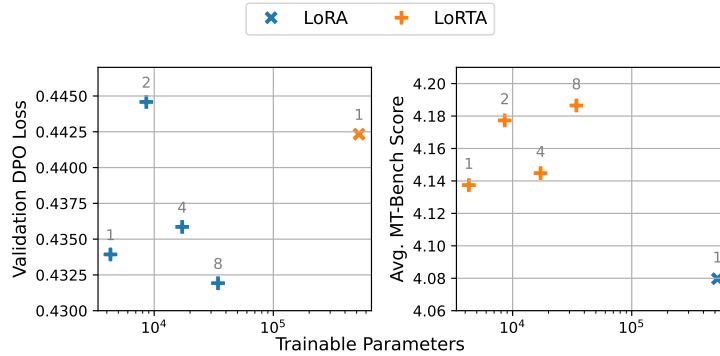


Figure 4. (Left) Mean DPO loss on held-out data from the orca dpo pairs dataset vs number of trainable parameters, lower is better. (Right) MT-Bench average scores Scores vs number of trainable parameters, higher is better.

LoRTA consistently outperformed LoRA across all ranks, including at rank 2 where it had shown higher DPO loss on the preference dataset. This improvement suggests enhanced out-of-distribution generalization capabilities for LoRTA adapters since MT-bench differs from the training dataset.

¹<https://huggingface.co/datasets/argilla/distilabel-intel-orca-dpo-pairs>

²<https://github.com/huggingface/trl>

4.4 PROTEIN FOLDING

Protein folding, the process by which a protein’s amino acid sequence determines its three-dimensional structure, is a fundamental problem in molecular biology. Accurate prediction of protein structures from their sequences has significant implications for understanding protein function and designing new proteins for therapeutic purposes. ESMFold (Lin et al., 2023) is a frontier model for this task trained in two stages. First, ESM-2, a BERT-based (Devlin et al., 2019) protein language model, is trained with the masked-language-modeling objective on amino acid sequences. This unsupervised pretraining allows the model to capture complex patterns and relationships within protein sequences. Remarkably, valuable structural information emerges in the model’s features during this process (Rao et al., 2020). In the second stage, ESM-2 is frozen, and a model head predicting three-dimensional protein structures is trained on top of language model features.

We re-train ESMFold in the second stage – fine-tuning ESM-2 parameters (we use ESM-2 35M instead of the ESM-2 3B model used in Lin et al. (2023) due to compute constraints) with LoRA and LoRTA instead of freezing them. We evaluate performance with the Local Distance Difference Test for $C\alpha$ atoms (LDDT- $C\alpha$) (Mariani et al., 2013) – that measures accuracy of predicted protein structures by comparing the distance between alpha carbons in predicted and true structures. LDDT- $C\alpha$ ranges from 0 (poor accuracy) to 1 (perfect match). See Appendix E.4 for experiment details.

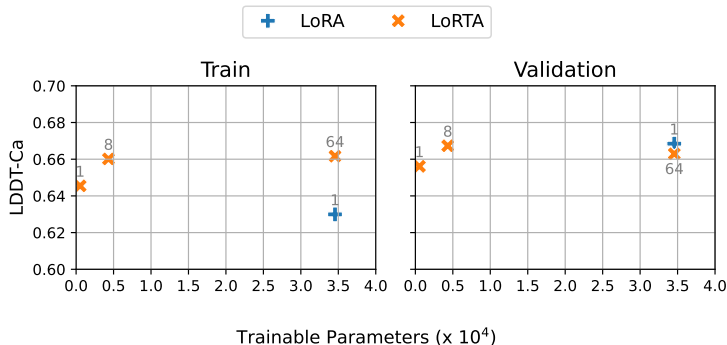


Figure 5. Mean LDDT- $C\alpha$ on train and held-out test sets. Higher is better. LoRTA rank 1 is competitive with LoRA rank 1 on the test set despite having 64x fewer parameters. Numbers on top of markers denote the adapter rank.

As shown in Figure 5, all tested LoRTA ranks outperform rank 1 LoRA on the training set; on the validation set, all tested LoRTA ranks are competitive with rank 1 LoRA. Notably, rank 1 LoRTA is competitive with rank 1 LoRA despite having an order of magnitude fewer parameters.

5 CONCLUSION

We have introduced LoRTA, a novel approach that employs a low-rank tensor model for LLM updates. By extending low-rank adaptation to higher-order tensors, LoRTA overcomes the inherent lower bounds on the number of trainable parameters while offering finer-grained control over adapter size. Our experiments across various benchmarks demonstrate that LoRTA achieves comparable and sometimes superior performance than baselines at a reduced parameter count.

Furthermore, we have shown that previous works have implicitly utilized low-rank tensor models with random factors. Nothing precludes our higher-order tensor model from using randomized factors for increased efficiency—a potential direction for future work that could further reduce computational overhead. Lastly, developing more efficient implementations of tensor operations that result in greater memory efficiency also remains a relevant future work direction which could make LoRTA even more suitable for resource-constrained environments.

REFERENCES

- Gustaf Ahlritz, Nazim Bouatta, Christina Floristean, Sachin Kadyan, Qinghui Xia, William Gerecke, Timothy J. O'Donnell, Daniel Berenberg, Ian Fisk, Niccolò Zanichelli, Bo Zhang, Arkadiusz Nowaczynski, Bei Wang, Marta M. Stepniewska-Dziubinska, Shang Zhang, Adegoke Ojewole, Murat Efe Guney, Stella Biderman, Andrew M. Watkins, Stephen Ra, Pablo Ribalta Lorenzo, Lucas Nivon, Brian Weitzner, Yih-En Andrew Ban, Shiyang Chen, Minjia Zhang, Conglong Li, Shuaiwen Leon Song, Yuxiong He, Peter K. Sorger, Emad Mostaque, Zhao Zhang, Richard Bonneau, and Mohammed AlQuraishi. Openfold: retraining alphafold2 yields new insights into its learning mechanisms and capacity for generalization. *Nature Methods*, 21(8):1514–1524, 2024. doi: 10.1038/s41592-024-02272-z. URL <https://doi.org/10.1038/s41592-024-02272-z>.
- Daniel Bershtatsky, Daria Cherniuk, Talgat Daulbaev, Aleksandr Mikhalev, and Ivan Oseledets. Lotr: Low tensor rank weight adaptation. *arXiv preprint arXiv:2402.01376*, 2024.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019. URL <https://arxiv.org/abs/1810.04805>.
- Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, et al. Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature Machine Intelligence*, 5(3):220–235, 2023.
- Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, Barret Zoph, Liam Fedus, Maarten P Bosma, Zongwei Zhou, Tao Wang, Emma Wang, Kellie Webster, Marie Pellat, Kevin Robinson, Kathleen Meier-Hellstern, Toju Duke, Lucas Dixon, Kun Zhang, Quoc Le, Yonghui Wu, Zhifeng Chen, and Claire Cui. GLaM: Efficient scaling of language models with mixture-of-experts. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 5547–5569. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/du22c.html>.
- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. Drop: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 2368–2378. Association for Computational Linguistics, 2019. doi: 10.18653/v1/N19-1246. URL <https://aclanthology.org/N19-1246/>.
- Raman Dutt, Linus Ericsson, Pedro Sanchez, Sotirios A Tsaftaris, and Timothy Hospedales. Parameter-efficient fine-tuning for medical image analysis: The missed opportunity. *arXiv preprint arXiv:2305.08252*, 2023.
- Ali Edalati, Marawan Gamal Abdel Hameed, and Ali Mosleh. Generalized kronecker-based adapters for parameter-efficient fine-tuning of vision transformers. In *2023 20th Conference on Robots and Vision (CRV)*, pp. 97–104, 2023. doi: 10.1109/CRV60082.2023.00020.
- William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39, 2022.
- Zihao Fu, Haoran Yang, Anthony Man-Cho So, Wai Lam, Lidong Bing, and Nigel Collier. On the effectiveness of parameter-efficient fine-tuning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 12799–12807, 2023.
- Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.

- Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. Minillm: Knowledge distillation of large language models. In The Twelfth International Conference on Learning Representations, 2024.
- Richard A Harshman and Margaret E Lundy. Parafac: Parallel factor analysis. Computational Statistics & Data Analysis, 18(1):39–72, 1994.
- Soufiane Hayou, Nikhil Ghosh, and Bin Yu. Lora+: Efficient low rank adaptation of large models. arXiv preprint arXiv:2402.12354, 2024.
- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a unified view of parameter-efficient transfer learning. arXiv preprint arXiv:2110.04366, 2021.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In International conference on machine learning, pp. 2790–2799. PMLR, 2019.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685, 2021.
- Shibo Jie and Zhi-Hong Deng. Fact: Factor-tuning for lightweight adaptation on vision transformer. Proceedings of the AAAI Conference on Artificial Intelligence, 37(1):1060–1068, Jun. 2023. doi: 10.1609/aaai.v37i1.25187. URL <https://ojs.aaai.org/index.php/AAAI/article/view/25187>.
- John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate protein structure prediction with alphafold. Nature, 596(7873):583–589, 2021. doi: 10.1038/s41586-021-03819-2. URL <https://doi.org/10.1038/s41586-021-03819-2>.
- Damjan Kalajdzievski. A rank stabilization scaling factor for fine-tuning with lora. arXiv preprint arXiv:2312.03732, 2023.
- Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke Hüllermeier. A survey of reinforcement learning from human feedback. arXiv preprint arXiv:2312.14925, 2023.
- Tamara G Kolda and Brett W Bader. Tensor decompositions and applications. SIAM review, 51(3): 455–500, 2009.
- Soroush Abbasi Koohpayegani, KL Navaneet, Parsa Nooralinejad, Soheil Kolouri, and Hamed Pirsiavash. Nola: Compressing lora using linear combination of random basis. In The Twelfth International Conference on Learning Representations, 2024.
- Dawid Jan Kopiczko, Tijmen Blankevoort, and Yuki Markus Asano. Vera: Vector-based random matrix adaptation. arXiv preprint arXiv:2310.11454, 2023.
- Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pp. 4582–4597, 2021a.
- Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pp. 4582–4597, 2021b.
- Baohao Liao, Shaomu Tan, and Christof Monz. Make pre-trained model reversible: From parameter to memory efficient fine-tuning. In Neural Information Processing Systems, 2023. URL <https://api.semanticscholar.org/CorpusID:258999607>.

- Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. Awq: Activation-aware weight quantization for on-device llm compression and acceleration. *Proceedings of Machine Learning and Systems*, 6: 87–100, 2024a.
- Xinyu Lin, Jinpeng Wang, Yong Zhang, Jingbo Zhang, and Ji-Rong Wen. Data-efficient fine-tuning for llm-based recommendation. *arXiv preprint arXiv:2401.17197*, 2024b.
- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, Allan dos Santos Costa, Maryam Fazal-Zarandi, Tom Sercu, Salvatore Candido, and Alexander Rives. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023. doi: 10.1126/science.ade2574. URL <https://www.science.org/doi/abs/10.1126/science.ade2574>.
- Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. *Advances in Neural Information Processing Systems*, 35:1950–1965, 2022.
- Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. Gpt understands, too. *AI Open*, 2023.
- Xinyin Ma, Gongfan Fang, and Xinchao Wang. Llm-pruner: On the structural pruning of large language models. *Advances in neural information processing systems*, 36:21702–21720, 2023.
- Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D. Lee, Danqi Chen, and Sanjeev Arora. Fine-tuning language models with just forward passes. *arXiv preprint arXiv:2305.17333*, 2023.
- Valerio Mariani, Marco Biasini, Alessandro Barbato, and Torsten Schwede. lddt: a local superposition-free score for comparing protein structures and models using distance difference tests. *Bioinformatics*, 29(21):2722–2728, 2013. doi: 10.1093/bioinformatics/btt473. URL <http://swissmodel.expasy.org/lddt>.
- Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, et al. Rwkv: Reinventing rnns for the transformer era. *arXiv preprint arXiv:2305.13048*, 2023.
- Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. Adapterfusion: Non-destructive task composition for transfer learning. *arXiv preprint arXiv:2005.00247*, 2020.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2383–2392. Association for Computational Linguistics, 2016. doi: 10.18653/v1/D16-1264. URL <https://aclanthology.org/D16-1264/>.
- Roshan Rao, Joshua Meier, Tom Sercu, Sergey Ovchinnikov, and Alexander Rives. Transformer protein language models are unsupervised structure learners. *bioRxiv*, 2020. doi: 10.1101/2020.12.15.422761. URL <https://www.biorxiv.org/content/early/2020/12/15/2020.12.15.422761>.
- Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In *Proceedings of the 2011 AAAI Spring Symposium Series*. Association for the Advancement of Artificial Intelligence, 2011.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.

- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pp. 234–241. Springer, 2015.
- Andreas Rücklé, Gregor Geigle, Max Glockner, Tilman Beck, Jonas Pfeiffer, Nils Reimers, and Iryna Gurevych. Adapterdrop: On the efficiency of adapters in transformers. *arXiv preprint arXiv:2010.11918*, 2020.
- Tara N Sainath, Brian Kingsbury, Vikas Sindhwani, Ebru Arisoy, and Bhuvana Ramabhadran. Low-rank matrix factorization for deep neural network training with high-dimensional output targets. In *2013 IEEE international conference on acoustics, speech and signal processing*, pp. 6655–6659. IEEE, 2013.
- Ying Sheng, Shiyi Cao, Dacheng Li, Coleman Hooper, Nicholas Lee, Shuo Yang, Christopher Chou, Banghua Zhu, Lianmin Zheng, Kurt Keutzer, et al. S-lora: Serving thousands of concurrent lora adapters. *arXiv preprint arXiv:2311.03285*, 2023.
- Nicholas D Sidiropoulos, Lieven De Lathauwer, Xiao Fu, Kejun Huang, Evangelos E Papalexakis, and Christos Faloutsos. Tensor decomposition for signal processing and machine learning. *IEEE Transactions on Signal Processing*, 65(13):3551–3582, 2017.
- Pablo Sprechmann, Alexander M Bronstein, and Guillermo Sapiro. Learning efficient sparse and low rank models. *IEEE transactions on pattern analysis and machine intelligence*, 37(9):1821–1833, 2015.
- Yang Sui, Miao Yin, Yu Gong, Jinqi Xiao, Huy Phan, and Bo Yuan. Elrt: Efficient low-rank training for compact convolutional neural networks. *arXiv preprint arXiv:2401.10341*, 2024.
- Simeng Sun, Dhawal Gupta, and Mohit Iyyer. Exploring the impact of low-rank adaptation on the performance, efficiency, and regularization of rlhf. *arXiv preprint arXiv:2309.09055*, 2023.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- John Thickstun. The transformer model in equations. *University of Washington: Seattle, WA, USA*, 2021.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Mojtaba Valipour, Mehdi Rezagholizadeh, Ivan Kobzyev, and Ali Ghodsi. Dylora: Parameter efficient tuning of pre-trained models using dynamic search-free low-rank adaptation. *arXiv preprint arXiv:2210.07558*, 2022.
- Gaël Varoquaux, Alexandra Sasha Luccioni, and Meredith Whittaker. Hype, sustainability, and the price of the bigger-is-better paradigm in ai, 2024. URL <https://arxiv.org/abs/2409.14160>.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*, 2018.
- Hongyi Wang, Saurabh Agarwal, Yoshiki Tanaka, Eric Xing, Dimitris Papailiopoulos, et al. Cuttlefish: Low-rank model training without all the tuning. *Proceedings of Machine Learning and Systems*, 5:578–605, 2023.
- Yifan Wang, Zixuan Xu, Yiren Shen, Yujun Zhu, and Hui Xiong. Insl: A data-efficient continual learning paradigm for fine-tuning large language models with instructions. *arXiv preprint arXiv:2403.11435*, 2024.

- Yeming Wen and Swarat Chaudhuri. Batched low-rank adaptation of foundation models. In The Twelfth International Conference on Learning Representations, 2024. URL <https://openreview.net/forum?id=w4ablTtZ2f>.
- Yifan Yang, Jiajun Zhou, Ngai Wong, and Zheng Zhang. Loretta: Low-rank economic tensor-train adaptation for ultra-low-parameter fine-tuning of large language models. arXiv preprint arXiv:2402.11417, 2024.
- Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. arXiv preprint arXiv:2106.10199, 2021.
- Longteng Zhang, Lin Zhang, Shaohuai Shi, Xiaowen Chu, and Bo Li. Lora-fa: Memory-efficient low-rank adaptation for large language models fine-tuning, 2023a.
- Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. In The Eleventh International Conference on Learning Representations, 2023b.
- Sheng Zhang, Xiaodong Liu, Jingjing Liu, Jianfeng Gao, Kevin Duh, and Benjamin Van Durme. Record: Bridging the gap between human and machine commonsense reading comprehension. arXiv preprint arXiv:1810.12885, 2018.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), Advances in Neural Information Processing Systems, volume 36, pp. 46595–46623. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/91f18a1287b398d378ef22505bf41832-Paper-Datasets_and_Benchmarks.pdf.
- Yitao Zhu, Zhenrong Shen, Zihao Zhao, Sheng Wang, Xin Wang, Xiangyu Zhao, Dinggang Shen, and Qian Wang. Melo: Low-rank adaptation is better than fine-tuning for medical image diagnosis. In 2024 IEEE International Symposium on Biomedical Imaging (ISBI), pp. 1–5. IEEE, 2024.

A TENSOR ALGEBRA

To facilitate our analysis, we briefly present some tensor algebra preliminaries and refer the reader to Sidiropoulos et al. (2017); Kolda & Bader (2009) for further details.

A N -order tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ is an N -way array indexed by i_1, i_2, \dots, i_N with elements $\mathcal{X}(i_1, i_2, \dots, i_N)$. It consists of N types of modes: $\mathcal{X}(:, i_2, \dots, i_N)$, $\mathcal{X}(i_1, :, \dots, i_N)$, \dots , $\mathcal{X}(i_1, i_2, \dots, :)$.

A rank-one tensor $\mathcal{Z} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ is the outer product of N vectors defined as:

$$\mathcal{Z} = \mathbf{a}_1 \circ \mathbf{a}_2 \circ \dots \circ \mathbf{a}_N, \quad (15)$$

where $\mathbf{a}_1 \in \mathbb{R}^{I_1}$, $\mathbf{a}_2 \in \mathbb{R}^{I_2}$, \dots , $\mathbf{a}_N \in \mathbb{R}^{I_N}$ and \circ denotes the outer product. The elementwise formula of the above expression is:

$$\mathcal{Z}(i_1, i_2, \dots, i_N) = \mathbf{a}_1(i_1) \mathbf{a}_2(i_2) \dots \mathbf{a}_N(i_N), \quad \text{for all } i_1, i_2, \dots, i_N, \quad (16)$$

Any tensor can be realized as a sum of N -way outer products (rank one tensors), i.e.

$$\mathcal{X} = \sum_{r=1}^R \mathbf{a}_1^r \circ \mathbf{a}_2^r \circ \dots \circ \mathbf{a}_N^r. \quad (17)$$

The above expression represents the *canonical polyadic decomposition* (CPD) or *parallel factor analysis* (PARAFAC) (Harshman & Lundy, 1994) of a tensor. The CPD elementwise representation is:

$$\mathcal{X}(i, j, k) = \sum_{r=1}^R \mathbf{A}_1(i_1, f) \mathbf{A}_2(i_2, f) \dots \mathbf{A}_N(i_N, f), \quad (18)$$

where $\mathbf{A}_n = [\mathbf{a}_n^1, \mathbf{a}_n^2, \dots, \mathbf{a}_n^F] \in \mathbb{R}^{I_n \times F}$, $n = 1, \dots, N$ are called the low rank factors of the tensor. A tensor can be fully characterized by its latent factors, so we can represent a tensor by its CPD model as:

$$\mathcal{X} = \llbracket \mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_N \rrbracket. \quad (19)$$

A tensor can be also represented as a set of matrices, by fixing all the modes but two as:

$$\mathcal{X}[:, :, i_3, \dots, i_N] = \mathbf{A}_1 (\text{Diag}(\mathbf{A}_3(i_3, :)) \odot \dots \odot \text{Diag}(\mathbf{A}_N(i_N, :))) \mathbf{A}_2^T, \quad (20)$$

where $\text{Diag}(\mathbf{A}_n(i_n, :))$ is the diagonal matrix with diagonal equal to $\mathbf{A}_n(i_n, :)$.

B ADDITIONAL RELATED WORK

Model Compression While these techniques differ from PEFT in that they focus on reducing the requirements of a trained model rather than efficient adaptation, they offer valuable insights for developing more efficient PEFT approaches. Pruning and quantization are key techniques for compressing neural networks, that have also been extensively applied to LLMs. Pruning removes less important weights, with some methods achieving high compression rates, e.g. (Ma et al., 2023). Quantization reduces weight precision, decreasing model size and also allowing more efficient operations (Lin et al., 2024a). Knowledge distillation is an alternative approach that involves transferring knowledge from a large “teacher” model to a smaller “student” model (Gu et al., 2024).

Low Rank Training. Exploiting low rank structure to improve efficiency during both training and inference in deep models has long been studied (Sainath et al., 2013), and also combined with sparsity (Sprechmann et al., 2015). Recent advancements include Cuttlefish (Wang et al., 2023) and ELRT (Sui et al., 2024).

Data efficient fine tuning. An alternative approach to reducing fine-tuning costs is to reduce the amount of data. In this direction, Few-shot and continual learning approaches have been shown to be effective in LLM fine-tuning tasks (Lin et al., 2024b; Wang et al., 2024).

Efficient Architectures Another relevant direction in resource usage is using more efficient model architectures. Mixture of Experts (MoE) technique, implemented in models like Switch Transformers (Fedus et al., 2022) and GLaM (Du et al., 2022), has shown promise in scaling model capacity

while maintaining computational efficiency by activating only relevant sub-models for given inputs. There is also relevant work on non-transformer architectures, such as RWKV (Peng et al., 2023) and Mamba (Gu & Dao, 2023), which combines the strengths of RNNs and Transformers to achieve efficient inference and training.

C OTHER TENSOR LOW RANK MODELS IN PEFT

D PARAMETER EFFICIENCY COMPARISON AGAINST LORA.

To fairly compare the parameter efficiency of LoRTA with LoRA, we adjust the tensor rank in LoRTA to match the effective total tensor rank in LoRA, which is $r' = r \times 4L$ due to LoRA applying a rank r update to each of the $4L$ matrices individually. For a given tensor rank, LoRTA reduces the number of parameters from scaling $8dLr$ in LoRA to $4L(d(1 + 1/h) + h + L + 4)r$ in LoRTA (usually $d \gg L$ and $d \gg h$), achieving substantial parameter savings without compromising expressive power. For example, this amounts to a 47.6% reduction in a LLaMA2 7B model.

We provide a breakdown of the parameter savings achieved by our proposed method, LoRTA, compared to LoRA, by parameterizing the weight updates using low-rank tensor decompositions at different granularities. The table below summarizes the dimensions of the update tensors, the number of update tensors used, and the corresponding parameter savings when the tensor rank r matches the tensor rank of LoRA rank r . The first row corresponds to LoRA.

Table 4: Update Tensor Modes, Parameters, and Savings

Added Modes	Update Tensor Dimensions	Number of Update Tensors	Parameter Savings
	$d \times d$	$4L$	0
Heads	$d \times \frac{d}{H} \times H$	$4L$	$1 - \frac{d(1 + \frac{1}{H}) + H}{2dr}$
Heads, QKVP	$d \times \frac{d}{H} \times H \times 4$	L	$1 - \frac{d(1 + \frac{1}{H}) + H + 4}{2dr}$
Heads, QKVP, Layers	$d \times \frac{d}{H} \times H \times 4 \times L$	1	$1 - \frac{d(1 + \frac{1}{H}) + H + 4 + L}{2dr}$

E EXPERIMENTAL DETAILS

In this appendix, we provide further details on the experiments presented in the main paper.

E.1 NLU

In our GLUE experiments we implemented our method using Huggingface’s PEFT, VeRA Kopiczko et al. (2023) and LoreTTA Yang et al. (2024) codebases. Hyperparameters for each of the three settings reported are detailed below.

Hyperparameter	Value
α	16
Learning Rate	[2E-3, 5E-4]
Scheduler	Constant
Optimizer	AdamW
Number of Epochs	20
Batch Size	[16, 32]
Warmup Steps	500

Table 5: Hyperparameter configurations for RoBERTa Base on the GLUE benchmark following the setup reported by Yang et al. (2024), where only the batch size and learning rate are tuned for each task, selecting between two values based on validation performance. All other hyperparameters match those reported by Yang et al. (2024).

Hyperparameter	Value
α	[0.5 1.0 2.0 8.0]
Learning Rate	[5e-4, 1e-3, 5e-3, 1e-2]
Scheduler	Linear
Optimizer	AdamW
Number of Epochs	20
Batch Size	32
Warmup Ratio	0.06

Table 6: Hyperparameter configurations for RoBERTa Base on the GLUE benchmark following Bershtsky et al. (2024). A grid-search to set the learning rate and scale parameter for each task is conducted across the specified values.

Hyperparameter	SST-2	MRPC	CoLA	QNLI	RTE	STS-B
Optimizer	AdamW					
Warmup Ratio	0.06					
LR Schedule	Linear					
Epochs	10	40	40	20	40	20
Learning Rate (Head)	6E-3	3E-3	6E-3	2E-4	2E-3	2E-3
Learning Rate (Encoder)	1E-2	1E-2	1E-2	1E-2	2E-2	2E-2
Batch Size	32					

Table 7: Hyperparameter configurations for RoBERTa large on the GLUE benchmark. All other hyperparameters are taken from Kopiczko et al. (2023).

E.2 INSTRUCTION TUNING

For instruction tuning experiments we utilized Lightning AI’s LitGPT codebase and training recipe. Hyperparameters are detailed below.

Parameter	Value
α	16
Learning Rate	0.01
Scheduler	Cosine
Optimizer	AdamW
Weight Decay	0.01
Number of Epochs	1
Steps	51000
Batch Size	16
Warmup Steps	318

Table 8: Hyperparameter configurations for LLama2-7B on the Alpaca dataset.

E.3 DPO

For preference optimization experiments we utilized using Huggingface trl library’s dpo implementation and example script. Hyperparameters are detailed below.

E.4 PROTEIN FOLDING

For protein folding experiments, we utilized OpenFold Ahdriz et al. (2024) training code and datasets. The following modifications were made to the ESMFold model architecture due to limited compute resources: a) utilize 12 Evoformer layers instead of the 48 used in (Lin et al., 2023) b) utilize ESM-2 35M instead of ESM-2 3B c) maintain outer product mean implementation from (Jumper et al., 2021).

Table 9: Hyperparameter configurations for LLama2-7B on intel orca DPO pairs.

Parameter	Value
α	16
Learning Rate	0.00005
Scheduler	Cosine
Optimizer	AdamW
Weight Decay	0
Number of Epochs	1
Batch Size	16
Warmup Steps	200

Optimizer and learning rate scheduler were identical to (Jumper et al., 2021). Models were trained for 850,000 steps with batch size of 32. Validation metrics were computed using the validation set from (Ahdritz et al., 2024).

Preliminary experiments revealed that higher values of α yield better results in this setting. α for LoRA and LoRTA experiments was then selected in multiple stages. Initially, models were trained with α values of $256 \times r$ and $128 \times r$, and the best-performing model was chosen. If both configurations diverged, α was halved, and models were retrained with the next lower pair (e.g., $64 \times r$ and $32 \times r$). This halving process continued until a convergent model was found. See Table 10 for the selected α values across experiments.

Table 10: Selected α and LDDT-CA for protein folding models.

Model	α	Validation LDDT- $C\alpha$
LoRA ($r = 1$)	128	0.668
LoRTA ($r = 64$)	128	0.663
LoRTA ($r = 8$)	256	0.667
LoRTA ($r = 1$)	2	0.656

F ADDITIONAL RESULTS

Figure 6 shows that Validation gains were primarily driven by reduced training error, though generalization slightly worsened, particularly at rank 2. On the other hand, as already mentioned, MT-bench performance was comparable or superior for LoRTA across all ranks, as shown in Figure 7.

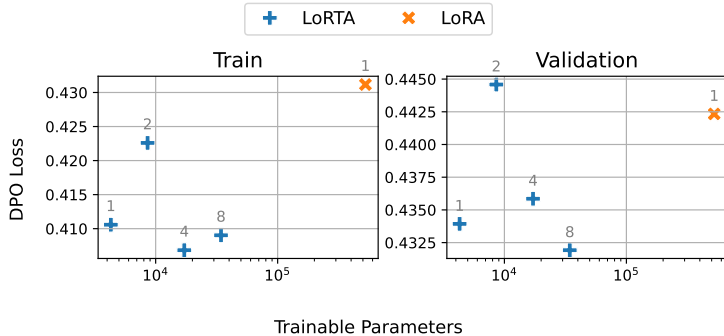


Figure 6. Mean DPO loss on the training (Left) and on held-out data (Right) from the orca dpo pairs dataset vs number of trainable parameters, lower is better.

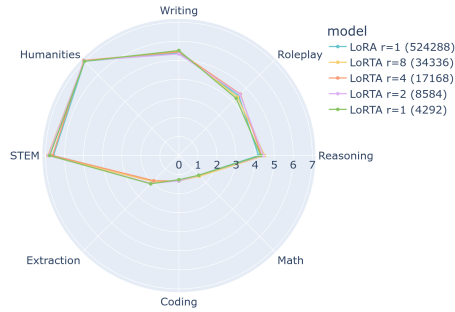


Figure 7. Performance on MT-Bench Zheng et al. (2023) for llama2-7b Touvron et al. (2023) models fine-tuned with LoRA and LoRTA using dpo on intel orca pairs. Average score per task. Higher is better.

G PRACTICAL IMPLICATIONS OF ADAPTER SIZE REDUCTION

The reduction in adapter size is primarily motivated by the need to improve task-switching efficiency and minimize storage requirements in scenarios involving a large number—potentially thousands—of adapters. Frequent CPU-GPU transfers for loading adapters in such settings can introduce significant overhead. By further compressing parameters, it becomes feasible for thousands of customized models to coexist with a base LLM in GPU memory, substantially enhancing scalability and performance in multi-task environments.

During training, the reduction in GPU memory usage from shrinking optimizer states is marginal for parameter reductions beyond LoRA. Memory consumption in these cases is dominated by activations and caches stored during forward and backpropagation. Additional memory savings could be achieved by compressing activations or gradients, leveraging the low-rank structure of updates, or dynamically recomputing them. While our model features fewer trainable parameters and could theoretically benefit from the efficient tensor CP structure, such as faster training and lower memory usage, these advantages are not yet realized due to the limitations of our current implementation. Future work will focus on optimizing this implementation. Future work will address these optimizations. However, the reduced parameter count already provides lower storage requirements and faster I/O.

We conducted hardware profiling to compare the performance of our LoRTA implementation against LoRA using HuggingFace PEFT. The results demonstrate negligible differences in resource consumption between the two methods. The slight gap in training time for LoRTA can be addressed through further optimizations, ranging from leveraging tools like Torch Compile, to implementing our CP tensor adapter model more efficiently.

Rank	Method	GPU Mem. (GB)	FLOPs (avg)	MACs (avg)	Time (s/step)
4	LoRA	12.84	272	136	0.07
	LoRTA	12.88	272	136	0.14
64	LoRA	13.08	276	138	0.09
	LoRTA	12.98	273	136	0.14

Table 11: Maximum GPU memory usage (GB), average FLOPs(GB), MACs(GB), and training time (seconds per step) for LoRA and LoRTA.