

# 1LoRA: Summation Compression for Very Low-Rank Adaptation

Alessio Quercia<sup>1,3</sup> Zhuo Cao<sup>1</sup> Arya Bangun<sup>1</sup> Richard D. Paul<sup>1,4</sup>  
 Abigail Morrison<sup>2,3</sup> Ira Assent<sup>1,5</sup> Hanno Scharr<sup>1</sup>

<sup>1</sup> IAS-8 <sup>2</sup> IAS-6, Forschungszentrum Juelich, Juelich, Germany

<sup>3</sup> Dept. of Computer Science, RWTH Aachen University, Aachen, Germany

<sup>4</sup> Dept. of Statistics, LMU Munich, Munich, Germany

<sup>5</sup> Dept. of Computer Science, Aarhus University, Aarhus, Denmark

{a.quercia, z.cao, a.bangun, r.paul, a.morrison, i.assent, h.scharr}@fz-juelich.de

## Abstract

Parameter-Efficient Fine-Tuning (PEFT) methods have transformed the approach to fine-tuning large models for downstream tasks by enabling the adjustment of significantly fewer parameters than those in the original model matrices. In this work, we study the "very low rank regime", where we fine-tune the lowest amount of parameters per linear layer for each considered PEFT method. We propose 1LoRA (Summation Low-Rank Adaptation), a compute, parameter and memory efficient fine-tuning method which uses the feature sum as fixed compression and a single trainable vector as decompression. Differently from state-of-the-art PEFT methods like LoRA, VeRA, and the recent MoRA, 1LoRA uses fewer parameters per layer, reducing the memory footprint and the computational cost. We extensively evaluate our method against state-of-the-art PEFT methods on multiple fine-tuning tasks, and show that our method not only outperforms them, but is also more parameter, memory and computationally efficient. Moreover, thanks to its memory efficiency, 1LoRA allows to fine-tune more evenly across layers, instead of focusing on specific ones (e.g. attention layers), improving performance further.

## 1. Introduction

In recent years, the rapid development of large models has transformed multiple fields in artificial intelligence, especially in natural language processing [1, 33] and computer vision [14, 27, 39]. However, adapting these models for specific tasks often demands significant computational power and extensive labeled data, making their deployment challenging for many applications. To mitigate these issues, Parameter Efficient Fine-Tuning (PEFT) methods [11, 12, 15, 20, 36, 41] have emerged as a promising approach to modify pre-trained models with minimal adjust-

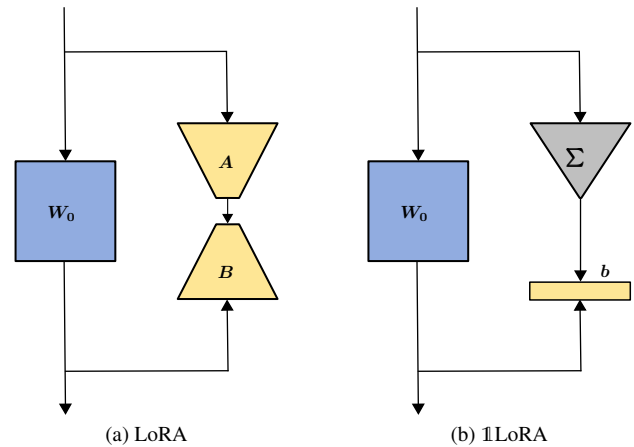


Figure 1. Comparing our method 1LoRA to LoRA. Left: LoRA learns the low-rank decomposition  $\Delta W = BA$ , where  $A \in \mathbb{R}^{r \times k}$  and  $B \in \mathbb{R}^{d \times r}$ . Right: 1LoRA replaces the matrices  $A$  and  $B$  with a sum over the input features  $x$  as compression and a learnable vector  $b \in \mathbb{R}^{1 \times d}$  as decompression:  $\Delta W = b\mathbb{1}^T$ , where the feature sum is  $\mathbb{1}^T x = \sum_{i=1}^k x_i$ , with  $\mathbb{1}$  being a vector of length  $k$  containing only ones. This reduces the trainable parameters from  $r \times k + d \times r$  in LoRA to  $d$  per layer for 1LoRA.

ments to their parameters, thereby reducing computational and memory demands while preserving performance.

Among various PEFT approaches, BitFit [41] has gained notable attention. It streamlines fine-tuning by updating only the bias parameters in transformer models, offering a lightweight alternative achieving competitive performance with minimal resource consumption. DiffFit [36] extends this method by adding a scaling factor to linear layers and also updating normalization layers, enabling rapid adaptation of large pre-trained diffusion models to new domains. These methods significantly speed up training and reduce model storage requirements, allowing practitioners to fine-tune models effectively without extensive retraining of all

parameters. Yet, because weight matrices are not changed they may not reach the performance achieved by other methods, such as Low-Rank Adaptation (LoRA) [11, 28].

LoRA [11] and its variants [6, 9, 12, 15, 19, 20] introduce low-rank matrices for model weight updates during fine-tuning, enabling efficient training with fewer parameters, see Figure 1 (left). This approach not only enhances the adaptability of large models but also reduces training time and mitigates overfitting. Compared to BitFit, LoRA achieves stronger performance in certain tasks and demonstrates robust generalization capabilities [31, 37]. This performance boost, however, requires more trainable parameters, which has been shown to cause training instability and slower convergence [5]. Additionally, Pu et al. [28] observe that LoRA is robust to parameter reduction for text classification and generation tasks, allowing for reductions beyond the original proposed ones. This raises the question: can we maintain or improve performance while reducing LoRA’s trainable parameters to levels comparable to BitFit?

In this work, we propose  $\mathbb{1}$ LoRA, a method centered on very low-rank adaptation, as depicted in Figure 1. We reduce the parameter count to match that of BitFit [41], making it significantly lower than that required by other low-rank adaptation methods.  $\mathbb{1}$ LoRA applies a predefined non-trainable function for input compression, *i.e.* the sum over all input features  $x$  and utilizes a single trainable vector per layer for decompression. As we demonstrate in our extensive, comparative experiments  $\mathbb{1}$ LoRA provides enhanced performance and memory efficiency in different tasks.

Our contributions can be summarized as follows:

- We propose  $\mathbb{1}$ LoRA, a very low-rank adaptation method which fixes the compression mechanism to a simple feature sum and a vector of  $d$  parameters for decompression, where  $d$  is the number of output features, significantly reducing the number of trainable parameters.
- Due to its low memory footprint at training time,  $\mathbb{1}$ LoRA enables fine-tuning across more layers in larger models (*e.g.*, LLaMA2 7 and 13 billion in our experiments), distributing adjustments across multiple layers instead of focusing solely on limited ones (*e.g.* the attention layers), thereby achieving enhanced performance.
- $\mathbb{1}$ LoRA adds no overhead at inference time, as the update  $\Delta W$  can be merged with the initial weight matrix  $W_0$  in the same way as *e.g.* done in LoRA [11].
- Our method is simple to implement, making it accessible for practical deployment.
- We conduct extensive empirical evaluations demonstrating that  $\mathbb{1}$ LoRA outperforms state-of-the-art methods in terms of effectiveness, and parameter, memory and computational efficiencies across multiple tasks and datasets.

## 2. Related Work

PEFT methods aim to adapt large pre-trained models for specific downstream tasks by injecting a small number of learnable parameters while freezing the majority of the model’s parameters. This approach addresses the resource inefficiency associated with full fine-tuning methods, where all model parameters are tuned for the new task.

Various strategies are employed to integrate learnable parameters in PEFT. For instance, Adapter [10] methods introduce small, fine-tunable modules across different layers of Transformer models, optimizing only these new modules for downstream adaptation. Prompt-tuning [17] and Prefix-tuning [18] involve inserting a set of trainable tokens into the input layer of Transformers, focusing training efforts on these tokens only. More recently, LoRA [11] reparameterizes certain pre-trained weights using low-rank matrices, tuning only these additional branches to adapt to new tasks. BitFit [41] proposes to fine-tune only bias terms, minimizing the number of trainable parameters. Although updating only a fraction of the parameters, these PEFT methods achieve comparable or even superior performance to full fine-tuning approaches.

While PEFT methods were initially developed within the natural language processing domain, recent works have extended their application to computer vision tasks. For example, [42] apply LoRA on the Segment Anything model [14], and DiffFit [36] transfers BitFit to fine-tune the diffusion model. Moreover, PEFT methods have recently gained interest for subspace Bayesian inference in large-scale Transformers [23, 26, 38], demonstrating the effectiveness of this efficient tuning paradigm across multiple domains.

## 3. Very Low-Rank Adaptation

In this section, we first analyze existing PEFT methods in regards to their fine-tuning approach and number of trainable parameters. We then motivate and describe technical details of our memory efficient very-low rank adaptation method,  $\mathbb{1}$ LoRA. Finally, we compare our method with other methods in terms of the number of parameters and computational and memory efficiency, respectively.

### 3.1. PEFT methods

We start by reviewing the technical details of state-of-the-art PEFT methods, analyzing in particular their computational and memory requirements. We provide an overview for easy comparison in Table 1.

**LoRA** [11] was one of the first methods proposed for parameter efficient fine-tuning. Instead of fine-tuning the entire parameter matrix  $W_0 \in \mathbb{R}^{d \times k}$ , LoRA fine-tunes an ad-

Table 1. Analysed PEFT methods: number of trainable parameters in the "very low-rank regime" with  $r = 1$  (LoRA, VeRA and DoRA) and  $\hat{r} = \lfloor \sqrt{d} \rfloor$  (MoRA).

Method	Equation	#Params $\times$ Linear
LoRA [11]	$W_0x + \underline{BA}x$	$k + d$
DoRA [20]	$W_0x + m \frac{W_0 + \underline{BA}}{\ W_0 + \underline{BA}\ _c} x$	$d + k + d$
VeRA [15]	$W_0x + \underline{bBd}Ax$	$1 + d$
MoRA <sub>1,6</sub> [12]	$W_0x + f_d(\underline{M}f_c(x))$	$\hat{r}^2$
BitFit [41]	$W_0x + \underline{\beta}$	$d$
DiffFit [36]	$\underline{\gamma}(W_0x + \underline{\beta}), \underline{\text{norms}}$	$> 2d$
$\mathbb{1}$ LoRA	$W_0x + \underline{b}\mathbb{1}^T x$	$d$
All	$\underline{W_0}x + \underline{\beta}, \underline{\text{norms}}$	$> k \times d + d$

ditive shift matrix  $\Delta W \in \mathbb{R}^{d \times k}$  as follows:

$$o = W_0x + \underline{\Delta W}x = (W_0 + \underline{\Delta W})x \quad (1)$$

where  $x \in \mathbb{R}^k$  and  $o \in \mathbb{R}^d$  are input and output features, respectively. We underline symbols indicating trainable parameters to distinguish them from pretrained, frozen parameters. LoRA decomposes the shift matrix  $\Delta W$  by two rectangular matrices:

$$\Delta W = \underline{BA}, \quad \text{with } A \in \mathbb{R}^{r \times k}, B \in \mathbb{R}^{d \times r}. \quad (2)$$

Matrix  $A$  reduces the rank of input from  $k$  to  $r$  and matrix  $B$  decompresses it to the size of output  $d$ .

**DoRA** (weight-decomposed low-rank adaptation) [20] builds on LoRA and proposes to decompose the pretrained matrices into magnitude  $m \in \mathbb{R}^d$  and direction  $(W_0 + \underline{BA}) / (\|W_0 + \underline{BA}\|_c)$ . Here, the norm is defined as vector-wise norm of a matrix across each column.

**VeRA** (Vector-based Random-matrix Adaptation) [15] builds on top of LoRA using fixed random compression and decompression matrices ( $A$  and  $B$ ), where the only trainable parameters are two vectors ( $d \in \mathbb{R}^r$  and  $b \in \mathbb{R}^d$ ), respectively after the compression and decompression.

**MoRA** [12], similar to VeRA, also uses fixed compression  $f_c$  and decompression  $f_d$  schemes, while only learning the inner low-rank squared matrix  $M \in \mathbb{R}^{\hat{r} \times \hat{r}}$ , where the rank is  $\hat{r} = \lfloor \sqrt{(k+d)r} \rfloor$ . The authors of [12] propose two versions, named types 1 and 6 in their Python code<sup>1</sup>. The version 'type 1 (Sharing)' is according to Eq. 6 in [12] and uses a sum within groups of features as compression and the decompression is a copy of learned features. The 'type 6

<sup>1</sup><https://github.com/kongds/MoRA/>

(RoPE based)' is a version according to Eq. 9 in [12] where the features are reshaped into a new axis as compression, processed, and then reshaped back. We name these methods MoRA<sub>1</sub> and MoRA<sub>6</sub>, respectively.

**BitFit** [41] is one of the first PEFT methods, proposing to fine-tune only the biases  $\beta \in \mathbb{R}^d$  of a pre-trained model.

**DiffFit** [36] extends BitFit to diffusion models by introducing a scaling vector  $\gamma \in \mathbb{R}^d$ . As a result, this approach adds an additional parameter that requires fine-tuning. Moreover, DiffFit trains the normalization layers and the class embedding, if present in the model.

In this paper, we explore the "very low-rank regime", where we use the smallest amount of parameters for each state-of-the-art method. We notice that BitFit [41] and VeRA [15] are the methods that use the least parameters in this scenario, with  $d$  and  $1+d$  per layer, when the rank is  $r = 1$ . All other methods use more parameters. VeRA smartly shares the constant random matrices  $A$  and  $B$  across modules with the same shape to reduce computation, however, it still uses considerable memory for matrix storage. BitFit, on the other hand, only fine-tunes the bias terms that remain constant across inputs and therefore lack generalization capabilities. Furthermore, we notice that MoRA [12] has memory and computational inefficiencies due to the grouping and duplication operations, whereas other state-of-the-art LoRA methods heavily rely on the decomposition matrices  $A$  and  $B$ , which requires them to use additional memory.

### 3.2. $\mathbb{1}$ LoRA

To reduce the memory consumption and to allow for better scaling to big models, we adopt a very low rank regime. Concretely, we propose to compress the input features to their sum  $\sum_{i=1}^k x_i$ . This sum can be expressed as  $\mathbb{1}^T x$ , where  $\mathbb{1}$  is the main diagonal of the all-positive quadrant. Our intuition is based on the fact that typical non-linearities (e.g. ReLU) move features towards this quadrant. Thus  $\mathbb{1}$  is a good candidate for a vector strongly correlating with the inputs when we have this prior knowledge, and it is as good as any other random vector when we do not.

Our approach can be viewed an extreme case of MoRA [12] which compresses the input features by computing the sum within specific groups of features and decompresses them by duplicating their outputs. Importantly, in contrast to MoRA,  $\mathbb{1}$ LoRA uses a single trainable vector  $b$  of  $d$  parameters for decompression, thereby retaining performance in fine-tuning. Using vector  $\mathbb{1}$  of length  $k$  containing only ones, the shift matrix can be written as

$$\Delta W = \underline{b}\mathbb{1}^T \quad \mathbb{1} \in \{1\}^{k \times 1}, b \in \mathbb{R}^{d \times 1} \quad (3)$$

Figure 1 depicts  $\mathbb{1}$ LoRA in comparison to LoRA [11]. In practice,  $\mathbb{1}$ LoRA only introduces a trainable vector of  $d$  pa-

rameters per layer. The compression can be implemented in PyTorch [25] as a sum over the features, removing the need for additional compression and decompression matrices, as in most state-of-the-art methods [11, 12, 15, 20]. This makes our method more computationally and memory efficient.

For fast and memory efficient inference, weights can be merged as the right side in Equation (1), *i.e.* the fine-tuned weight matrix  $W_{ft}$  is given as  $W_{ft} := W_0 + \Delta W$ . In our experiments, this efficient implementation has been applied to the competitors LoRA, MoRA, and DoRA [11, 12, 20].

### 3.3. 1LoRA Method Analysis

Our method relates to LoRA [11] with rank  $r = 1$  which fine-tunes the low-rank matrices A and B. Compared to LoRA, 1LoRA uses  $k$  fewer parameters per layer (see Table 1), while removing the need for the matrices A and B. MoRA (Eq. 6, [12]) introduces constant, non-trainable compression and decompression functions: a sum over feature groups for compression and a duplication of learned features for decompression. Our approach goes to the most extreme case, grouping features by summing across a single group that includes all features. Simultaneously, it eliminates the need for copying learned features by introducing the vector  $b$ , which scales the feature sum into  $d$  parameters during decompression. Although MoRA can be made comparable to 1LoRA in terms of parameters by using rank  $\hat{r} = \lfloor \sqrt{d} \rfloor$ , 1LoRA is more computationally efficient as it avoids the need for grouping features, learning the inner features, and copying them back to match the required shape.

Among low-rank adaptation methods, only VeRA has a comparable number of parameters. However, it uses random projection matrices, making the method more memory-hungry, initialization-dependent, and less interpretable. By contrast, using the feature sum as input for the trainable decompression, 1LoRA suggests that it suffices to learn to shift the pre-trained parameters by scaled feature sum of samples drawn from a dataset to learn the new task.

In comparison to BitFit, which only fine-tunes the biases  $\beta$ , 1LoRA can better learn from the data as it uses the sum of the features, resulting in better performance, while slightly sacrificing training speed, see experiments in Section 4. Lastly, our method is complementary to state-of-the-art methods that fine-tune additional terms, like biases (BitFit [41]) or normalization layers (DiffFit, [36]), and to methods building on top of LoRA [6, 9, 11, 19, 20].

## 4. Empirical evaluation

We benchmark 1LoRA against state-of-the-art PEFT methods in the "very low-rank regime" on multiple fine-tuning tasks and models: Monocular Depth Estimation with

Table 2. RMSE ( $\downarrow$ ) and AbsRel ( $\downarrow$ ) of DepthAnything model pre-trained on KITTI and fine-tuned to NYU using PEFT methods. LoRA, VeRA and DoRA with rank  $r = 1$ , MoRA with rank  $\hat{r} = \lfloor \sqrt{d} \rfloor$  and  $\tilde{d} = \hat{r}^2$  parameters.

Method	KITTI $\rightarrow$ NYU		NYU $\rightarrow$ KITTI	
	RMSE ( $\downarrow$ )	AbsRel ( $\downarrow$ )	RMSE ( $\downarrow$ )	AbsRel ( $\downarrow$ )
MoRA <sub>1</sub>	0.248 $\pm$ 0.00066	0.0670 $\pm$ 0.00018	2.509 $\pm$ 0.00455	0.0587 $\pm$ 0.00011
MoRA <sub>6</sub>	0.277 $\pm$ 0.00021	0.0774 $\pm$ 0.00009	2.971 $\pm$ 0.00239	0.0739 $\pm$ 0.00009
LoRA	0.250 $\pm$ 0.00087	0.0678 $\pm$ 0.00028	2.353 $\pm$ 0.00315	0.0583 $\pm$ 0.00010
VeRA	0.244 $\pm$ 0.00059	0.0668 $\pm$ 0.00013	2.304 $\pm$ 0.00524	0.0577 $\pm$ 0.00013
DoRA	0.249 $\pm$ 0.00088	0.0675 $\pm$ 0.00031	2.345 $\pm$ 0.00275	0.0581 $\pm$ 0.00004
BitFit	0.349 $\pm$ 0.00073	0.1015 $\pm$ 0.00022	4.806 $\pm$ 0.00392	0.1271 $\pm$ 0.00027
DiffFit	0.314 $\pm$ 0.00078	0.0895 $\pm$ 0.00030	3.904 $\pm$ 0.00539	0.0988 $\pm$ 0.00024
<b>1LoRA</b>	<b>0.238<math>\pm</math>0.00050</b>	<b>0.0647<math>\pm</math>0.00010</b>	<b>2.203<math>\pm</math>0.00780</b>	<b>0.0533<math>\pm</math>0.00025</b>
All	0.210 $\pm$ 0.00064	0.0571 $\pm$ 0.00033	1.916 $\pm$ 0.00303	0.0463 $\pm$ 0.00010

DepthAnything [39], Image Classification with ViT-Base [7], Mathematical Reasoning with LLaMA2 7b and 13b [34], and Image Generation with DiT-XL-2-256x256 [27]. For each method we use the lowest amount of parameters per layer. We compare 1LoRA with LoRA ( $r = 1$ ), VeRA ( $r = 1$ ), DoRA ( $r = 1$ ), MoRA ( $\hat{r} = \lfloor \sqrt{d} \rfloor$ ) types 1 and 6 (see Section 3.1), BitFit and DiffFit. We report results in forms of tables and bubble plots, where the mean of 4 different independent runs is reported for each method. In each plot, we report the best mean evaluation metrics as y-axis, and their required total (training and validation) wall clock time as x-axis, GPU (A100 40GB) memory consumption as bar plot, and number of trainable parameters as bubble area. Note that parameters are only a proxy for memory consumption, which can be higher depending on computational graph complexity. We report metrics' standard deviations in the table and plots (as error bars).

### 4.1. Monocular Depth Estimation

As first experiment, we investigate the Monocular Depth Estimation (MDE) fine-tuning task. Here, we start from a pre-trained metric depth model and fine-tune it to a new dataset without introducing new untrained modules like adapters. We use the model and training settings of Depth Anything [39], which are based on DINOv2 [24], MiDaS [3, 29], DPT [30] and ZoeDepth [2]. We adopt widely used MDE benchmark datasets: NYU (indoor) [22], KITTI (driving) [8] and DIODE (outdoor) [35]. We test the following fine-tuning cases: from NYU to KITTI, from KITTI to NYU, from NYU to DIODE Outdoor, from KITTI to DIODE Outdoor. Table 2 shows the Root Mean Squared Error (RMSE) and Absolute Relative Error (AbsRel) for the first two cases, respectively. Here, 1LoRA outperforms all the competitors and is the closest to fine-tuning the entire model. Additionally, Figure 2a and Figure 2b show the training time, number of parameters and GPU memory needed to achieve the minimum Root Mean Square Error

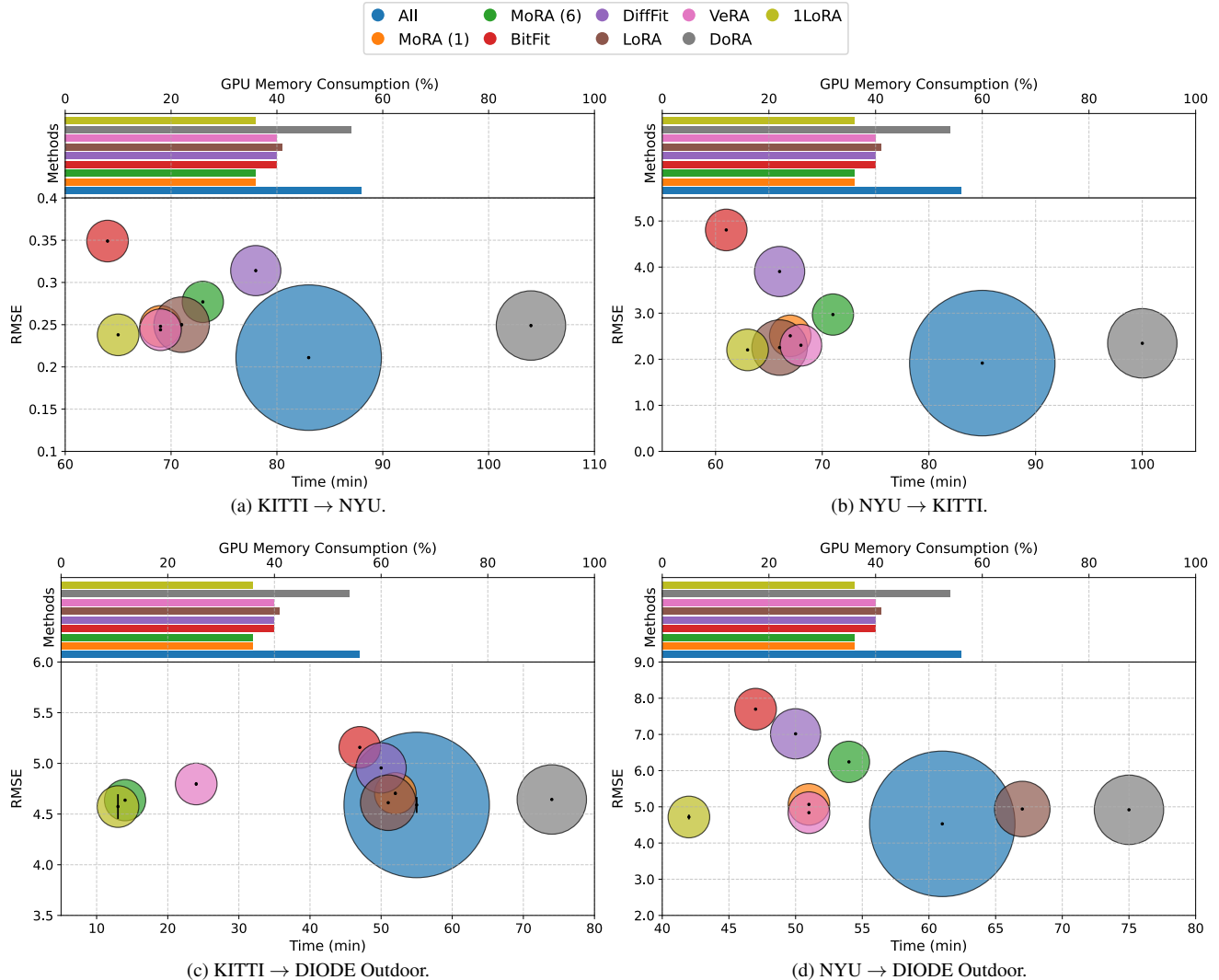


Figure 2. RMSE ( $\downarrow$ ) of pre-trained DepthAnything model fine-tuned using PEFT methods. Bubble size is proportional to the number of parameters, except for "All", which is capped due to space limitations. Bottom left (and smallest bubble) is better.

(RMSE) for each method in the first two cases. From the figure, it emerges that not only is our method  $\bullet$  the closest to full fine-tuning  $\bullet$ , but it is also (1) faster than all competitors but BitFit  $\bullet$ , (2) it is the method with the lowest memory consumption together with MoRA  $\bullet$ , (3) it is the method fine-tuning the least parameters together with BitFit, MoRA and VeRA  $\bullet$ .

Figure 2c and Figure 2d show results for the fine-tuning cases from KITTI to DIODE Outdoor and from NYU to DIODE Outdoor, respectively. The figures confirm that our method  $\bullet$  outperforms the competitors in terms of efficacy and computational efficiency, in terms of memory usage it is the best on par with MoRA  $\bullet$ , and in terms of parameters it is again the best on par with BitFit  $\bullet$ , MoRA and VeRA  $\bullet$ .

## 4.2. Mathematical Reasoning

We study the Mathematical Reasoning task, fine-tuning 7 billion and 13 billion parameters pre-trained LLaMA2 [34] models to Meta-Math [40]. We use the fine-tuning settings of MoRA [12], where for the 7 billion model we use 16 GPUs with 4 images each and a total batch size of 128 with gradient accumulation, and for the 13 billion model we use 64 GPUs with 2 images each. In addition, we split the dataset into 80% for training and 20% for validation.

Figure 3 shows that 1LoRA  $\bullet$  outperforms all methods, except LoRA  $\bullet$ , in terms of validation loss, while using competitive computational time, parameter and memory budgets. LoRA achieves a slightly lower validation loss, while using more than 90% of the GPU memory compared

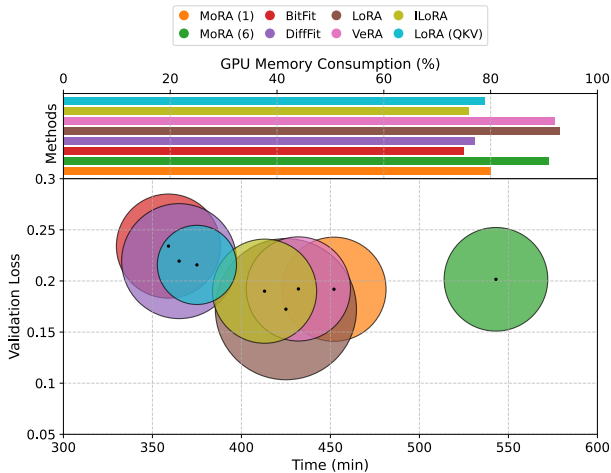


Figure 3. Validation loss of pre-trained LLaMA-2 7b fine-tuned to Meta-Math. Bubble size is proportional to the number of parameters. Bottom left (and smallest bubble) is better.

to less than 80% of 1LoRA. On the other hand, when limiting the memory budget of LoRA by applying it only to the attention layers, here named LoRA (QKV) ●, 1LoRA results in a better performance while still using slightly less GPU memory, showing that 1LoRA is more memory efficient than LoRA. Moreover, our method is only slower than BitFit ●, DiffFit ● and LoRA (QKV), while being better in terms of quality. In terms of trainable parameters, 1LoRA is again comparable to BitFit and MoRA ●, and only second to LoRA (QKV), which is applied to fewer layers. Lastly, in terms of memory, 1LoRA is the second most efficient method, almost on par with the first one, BitFit.

Given that 1LoRA is more memory efficient than LoRA and LoRA (QKV), 1LoRA can be applied to more layers and modules than LoRA in bigger models, allowing a more fine-grained fine-tuning than LoRA, which needs to be limited to specific layers. Figure 4 exemplifies this 1LoRA benefit for fine-tuning a LLaMA-2 13b. LoRA cannot be applied to all linear layers as it would require too much GPU memory, therefore it has to be limited to specific layers (Q, K, V in this case). 1LoRA can be applied to all linear layers with a similar memory budget, while enabling a better performance. In fact, here, 1LoRA outperforms all applicable competitors while consuming a comparable memory budget and only being slightly slower. We report visual comparison of 1LoRA and LoRA QKV in Appendix C.

### 4.3. Image Generation

We investigate the use of PEFT methods for diffusion model fine-tuning for image-generation. We fine-tune an ImageNet-21k [32] pre-trained DiT-XL-2-256x256 [27] to

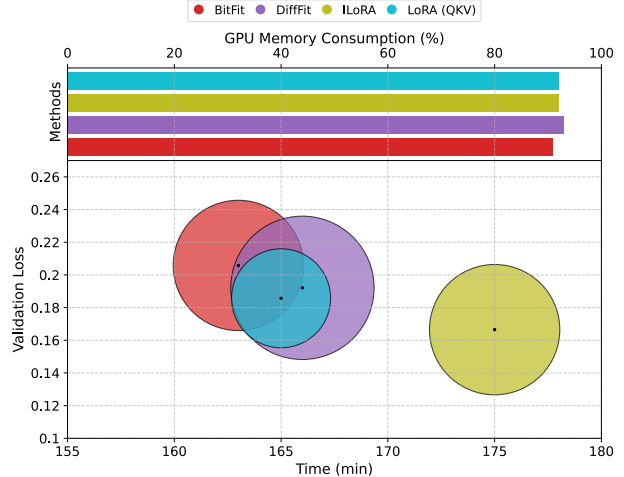


Figure 4. Validation loss of pre-trained LLaMA-2 13b fine-tuned to Meta-Math. Missing competitors cannot be fine-tuned with the given memory budget. Bubble size is proportional to the number of parameters. Bottom left (and smallest bubble) is better.

Food-101 [4] using the fine-tuning settings of DiffFit ● [36]. We run all experiments with a 24h time budget and we report results of a single run per method, given the high computational requirements. Figure 5 shows that 1LoRA ● outperforms all competitors in terms of FID, while being the most memory efficient method. Notice that 1LoRA also converges faster than LoRA ● and MoRA (6) ●. In addition, as for previous experiments, 1LoRA has the lowest number of parameters together with BitFit ●, MoRA ● and VeRA ●. Lastly, DoRA ● cannot be fine-tuned with the same fine-tuning settings as it goes out of memory.

### 4.4. Image Classification

We compare 1LoRA and competitor methods on Image Classification. We use a ViT-Base [7] model pre-trained on ImageNet-21k [32] and we fine-tune it to CIFAR10 and CIFAR100 [16], respectively, using AdamW [21] with learning rate 0.00002, weight decay 0.01, batch size 10, random resized crop to 224 and random horizontal flip.

Figure 6 and Figure 7 show the results for CIFAR10 and CIFAR100, respectively. From the figures, it emerges that 1LoRA ● is the second fastest method after BitFit ●, while being slightly better in terms of accuracy. LoRA ● and DoRA ● are the methods achieving the closest accuracy to fine-tuning the entire model, however they are also the slowest methods, even slower than full fine-tuning. In this experiment, DiffFit ● is slightly better than 1LoRA ●, however it is also slightly slower and it uses more parameters and memory. It is worth noting that 1LoRA is complementary to BitFit ● and DiffFit ●. They can be combined in

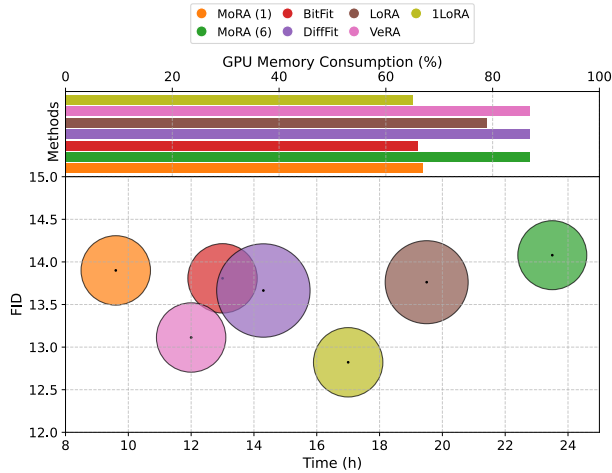


Figure 5. FID of pretrained DiT fine-tuned to Food-101. Note that DoRA is missing as it cannot be fine-tuned with the given memory budget. Bubble size is proportional to the number of parameters. Bottom left (and smallest bubble) is better.

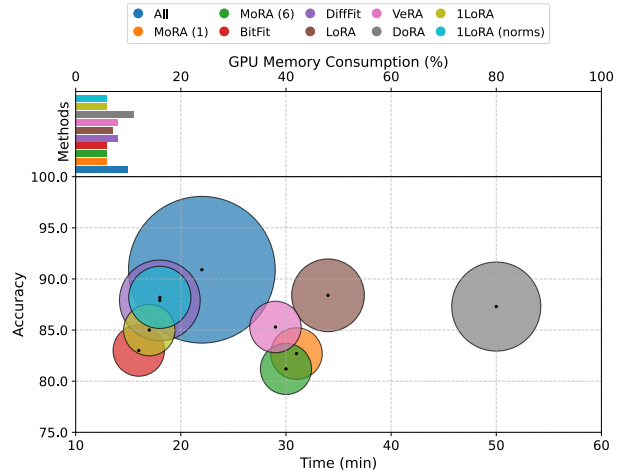


Figure 7. Accuracy ( $\uparrow$ ) of ViT-Base model pre-trained on ImageNet-21k and fine-tuned to CIFAR100. Bubble size is proportional to the number of parameters, except for "All", which is capped. Top left (and smallest bubble) is better.

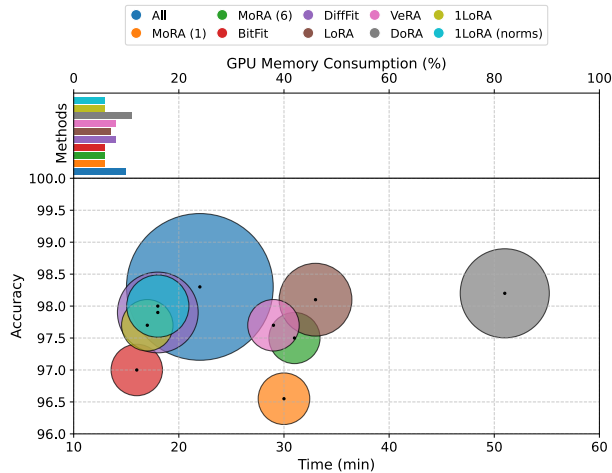


Figure 6. Accuracy ( $\uparrow$ ) of ViT-Base model pre-trained on ImageNet-21k and fine-tuned to CIFAR10. Bubble size is proportional to the number of parameters, except for "All", which is capped. Top left (and smallest bubble) is better.

many ways, as we study in Section 5.1. Here, the figures additionally report the most fruitful combination, 1LoRA (norms)  $\bullet$ , *i.e.* 1LoRA and unfrozen normalization layers. This combination is slightly better, being more parameter and memory efficient than DiffFit. In Section 5.1, we study further combinations. Lastly, in order to show that fine-tuning the backbone leads to significant advantages, we report accuracies for the fine-tuning of the classification head only, on CIFAR10 (90%) and CIFAR100 (66%).

## 5. Ablations

### 5.1. Complementary methods

We investigate all possible combinations of BitFit, DiffFit and 1LoRA, including individual components (the biases  $\beta$ , the scaling factors  $\gamma$ , and the norms, see Table 1), in order to understand which modules impact fine-tuning most. We run experiments for Image Classification, fine-tuning a ViT-Base from ImageNet-21k [32] to CIFAR10. We adopt the training configurations described in previous sections.

Figure 8 reports results for the Image Classification task on CIFAR10, where there is a clear advantage in preferring fine-tuning of normalization layers  $\bullet$  over biases  $\bullet$  or scaling factors  $\bullet$ . Combining these modules with 1LoRA leads to the best results, comparable to DiffFit  $\bullet$ , but with fewer parameters, less memory consumption and smaller computational cost. Other combinations do not seem to bring additional performance gains. We report additional results for this ablation on the MDE task in Appendix B.1.

### 5.2. PCA analysis

We performed PCA on the weight updates for our Monocular Depth Estimation and Classification experiments. We compare principal components (PCs) of a full-rank update against a random, summation (1LoRA) and learned (LoRA,  $r = 1$ ) compression vector and present results in Figure 9. For MDE, alignment of all compression vectors in terms of absolute cosine similarity to the PCs is mostly low, except for the second MLP layers. For classification, a similar pattern is visible except that the learned compression achieves

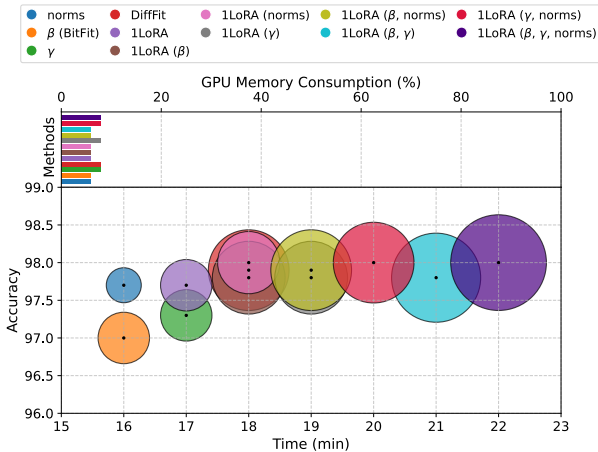


Figure 8. Accuracy ( $\uparrow$ ) of ViT-Base model pre-trained on ImageNet-21k and fine-tuned to CIFAR10, using all combinations of 1LoRA, BitFit and DiffFit, including individual components (biases  $\beta$ , scaling factors  $\gamma$  and normalization layers, *norms*). Bubble size is proportional to the number of parameters. Top left (and smallest bubble) is better.

notably higher similarity, especially also for the first MLP layer. Interestingly, the second MLP layer of a ViT block is the only layer in the architecture, where features went through a GELU activation function right before-hand. Effectively, the GELU activation concentrates the incoming features towards the positive “quadrant”, thus summation (*i.e.* the  $\mathbb{1}$  vector) is unlikely to be orthogonal to those incoming features, and orthogonality may be an issue to performance as it suppresses the incoming feature.

As to why 1LoRA yields good performance, we argue that summation seems to be a good a priori guess for the PCs of linear layers that are fed with GELU or ReLU activations. For the other layers, a fixed summation compression performs similarly to a random one. In Appendix B.2 we compare results for 1LoRA with fixed summation and random compression, showing that the former (*i.e.* 1LoRA) is at least as good as a random compression, while being faster. On the other hand, a learned compression only leads to slightly better alignment with the PCs (*cf.* learned compression in Figure 9), while requiring double the amount of parameters and more memory.

In conclusion, given the similarity in compression quality, using 1LoRA is preferable for memory and computational efficiencies.

## 6. Conclusion

PEFT methods have transformed the approach to fine-tuning large models for downstream tasks by enabling the

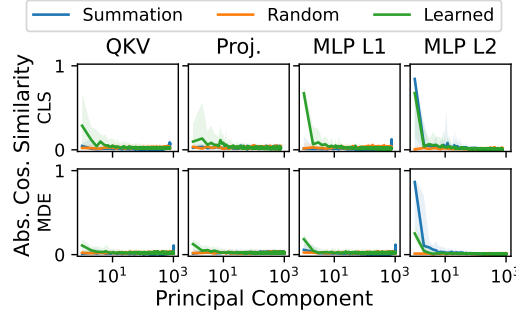


Figure 9. Alignment between PCs (sorted by singular values) of a full weight update and our candidate compression vectors.

adjustment of significantly fewer parameters than those in the original model matrices. In this work, we study the “very low rank regime”, where we fine-tune the lowest amount of parameters per layer for each considered PEFT method. Inspired by LoRA, VeRA, and the recent MoRA, we propose 1LoRA, a compute, parameter and memory efficient fine-tuning method which uses the feature sum as compression, and a trainable vector as decompression. Differently from state-of-the-art PEFT methods, 1LoRA uses only  $d$  parameters per layer, reducing the memory footprint and the computational cost.

Moreover, by using the feature sum, instead of random projection, as input for the trainable decompression, 1LoRA makes one step forward in interpretable fine-tuning, suggesting that learning to shift the pre-trained parameters by scaled feature sum coming from samples drawn from a new dataset is enough to learn the new task.

We extensively evaluate our method against state-of-the-art PEFT methods on multiple fine-tuning tasks and models, namely Monocular Depth Estimation with DepthAnything, Image Classification with ViT-Base, Mathematical Reasoning with LLaMA2 7b and 13b, and Image Generation with DiT-XL-2-256x256. We observe that our method not only outperforms state-of-the-art methods, but is also more parameter, memory and computationally efficient.

By reducing the number of trainable parameters per layer to  $d$ , 1LoRA enables the fine-tuning of big models (*e.g.* LLMs with billions of parameters), where other methods cannot be applied due to their larger memory needs. Additionally, we show that our method can be applied to more layers on bigger models, outperforming other competitors while using similar memory budget.

With this paper, we demonstrate the benefit of the very low-rank regime, where the lowest amount of parameters is fine-tuned, while maximizing the performance.



**Acknowledgements** A.Q. and R.D.P. were funded by the Helmholtz School for Data Science in Life, Earth, and Energy (HDS-LEE). The authors gratefully acknowledge the Gauss Centre for Supercomputing e.V. ([www.gauss-centre.eu](http://www.gauss-centre.eu)) for funding this project by providing computing time through the John von Neumann Institute for Computing (NIC) on the GCS Supercomputer JUWELS [13] at Jülich Supercomputing Centre (JSC).

## References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 1
- [2] Shariq Farooq Bhat, Reiner Birkel, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*, 2023. 4
- [3] Reiner Birkel, Diana Wofk, and Matthias Müller. Midas v3. 1—a model zoo for robust monocular relative depth estimation. *arXiv preprint arXiv:2307.14460*, 2023. 4
- [4] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative components with random forests. In *European Conference on Computer Vision*, 2014. 6
- [5] Guanzheng Chen, Fangyu Liu, Zaiqiao Meng, and Shangsong Liang. Revisiting parameter-efficient tuning: Are we really there yet? *arXiv preprint arXiv:2202.07962*, 2022. 2
- [6] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024. 2, 4
- [7] Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 4, 6
- [8] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *International Journal of Robotics Research (IJRR)*, 2013. 4
- [9] Han Guo, Philip Greengard, Eric P Xing, and Yoon Kim. Lq-lora: Low-rank plus quantized matrix decomposition for efficient language model finetuning. *arXiv preprint arXiv:2311.12023*, 2023. 2, 4
- [10] 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*, pages 2790–2799. PMLR, 2019. 2
- [11] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 1, 2, 3, 4
- [12] Ting Jiang, Shaohan Huang, Shengyue Luo, Zihan Zhang, Haizhen Huang, Furu Wei, Weiwei Deng, Feng Sun, Qi Zhang, Deqing Wang, et al. Mora: High-rank updating for parameter-efficient fine-tuning. *arXiv preprint arXiv:2405.12130*, 2024. 1, 2, 3, 4, 5
- [13] Stefan Kesselheim, Andreas Herten, Kai Krajsek, Jan Ebert, Jenia Jitsev, Mehdi Cherti, Michael Langguth, Bing Gong, Scarlet Stadler, Amirpasha Mozaffari, et al. Juwels booster—a supercomputer for large-scale ai research. In *International Conference on High Performance Computing*, pages 453–468. Springer, 2021. 9
- [14] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023. 1, 2
- [15] Dawid Jan Kopiczko, Tijmen Blankevoort, and Yuki Markus Asano. Vera: Vector-based random matrix adaptation. *arXiv preprint arXiv:2310.11454*, 2023. 1, 2, 3, 4
- [16] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *Citeseer*, 2009. 6
- [17] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021. 2
- [18] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*, 2021. 2
- [19] Yang Li, Shaobo Han, and Shihao Ji. Vb-lora: Extreme parameter efficient fine-tuning with vector banks. *arXiv preprint arXiv:2405.15179*, 2024. 2, 4
- [20] Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. Dora: Weight-decomposed low-rank adaptation. *arXiv preprint arXiv:2402.09353*, 2024. 1, 2, 3, 4
- [21] I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 6
- [22] Pushmeet Kohli Nathan Silberman, Derek Hoiem and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *ECCV*, 2012. 4
- [23] Emre Onal, Klemens Flöge, Emma Caldwell, Arsen Sheverdin, and Vincent Fortuin. Gaussian stochastic weight averaging for bayesian low-rank adaptation of large language models. In *Sixth Symposium on Advances in Approximate Bayesian Inference - Non Archival Track*, 2024. 2
- [24] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 4
- [25] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NIPS 2017 Autodiff Workshop*, 2017. 4
- [26] Richard D. Paul, Alessio Quercia, Vincent Fortuin, Katharina Nöh, and Hanno Schar. Parameter-efficient bayesian neural networks for uncertainty-aware depth estimation, 2024. 2
- [27] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023. 1, 4, 6

- [28] George Pu, Anirudh Jain, Jihan Yin, and Russell Kaplan. Empirical analysis of the strengths and weaknesses of peft techniques for llms. *arXiv preprint arXiv:2304.14999*, 2023. [2](#)
- [29] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE transactions on pattern analysis and machine intelligence*, 44(3):1623–1637, 2020. [4](#)
- [30] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 12179–12188, 2021. [4](#)
- [31] Olesya Razuvayevskaya, Ben Wu, João A Leite, Freddy Heppell, Ivan Srba, Carolina Scarton, Kalina Bontcheva, and Xingyi Song. Comparison between parameter-efficient techniques and full fine-tuning: A case study on multilingual news article classification. *Plos one*, 19(5):e0301738, 2024. [2](#)
- [32] Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lih Zelnik-Manor. Imagenet-21k pretraining for the masses, 2021. [6](#), [7](#)
- [33] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. [1](#)
- [34] 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. [4](#), [5](#)
- [35] Igor Vasiljevic, Nick Kolkin, Shanyi Zhang, Ruotian Luo, Haochen Wang, Falcon Z. Dai, Andrea F. Daniele, Mohamadreza Mostajabi, Steven Basart, Matthew R. Walter, and Gregory Shakhnarovich. DIODE: A Dense Indoor and Outdoor DDepth Dataset. *arXiv e-prints*, art. arXiv:1908.00463, 2019. [4](#)
- [36] Enze Xie, Lewei Yao, Han Shi, Zhili Liu, Daquan Zhou, Zhaoqiang Liu, Jiawei Li, and Zhenguo Li. DiffFit: Unlocking transferability of large diffusion models via simple parameter-efficient fine-tuning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4230–4239, 2023. [1](#), [2](#), [3](#), [4](#), [6](#)
- [37] Chunlei Xin, Yaojie Lu, Hongyu Lin, Shuheng Zhou, Huijia Zhu, Weiqiang Wang, Zhongyi Liu, Xianpei Han, and Le Sun. Beyond full fine-tuning: Harnessing the power of LoRA for multi-task instruction tuning. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 2307–2317, Torino, Italia, 2024. ELRA and ICCL. [2](#)
- [38] Adam X. Yang, Maxime Robeyns, Xi Wang, and Laurence Aitchison. Bayesian low-rank adaptation for large language models. In *The Twelfth International Conference on Learning Representations*, 2024. [2](#)
- [39] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. *arXiv:2401.10891*, 2024. [1](#), [4](#)
- [40] Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*, 2023. [5](#)
- [41] 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. [1](#), [2](#), [3](#), [4](#)
- [42] Zihan Zhong, Zhiqiang Tang, Tong He, Haoyang Fang, and Chun Yuan. Convolution meets lora: Parameter efficient finetuning for segment anything model. *arXiv preprint arXiv:2401.17868*, 2024. [2](#)

# Supplemental Material

## 1LoRA: Summation Compression for Very Low-Rank Adaptation

### A. Algorithm

As stated in the main paper, 1LoRA is straightforward to implement. We illustrate the python implementation in Listing 1. In our work, we applied the 1LoRA module to all linear layers in the model, excluding the classification layers, when present.

```

1 import torch
2 import torch.nn as nn
3
4 class ILoRA(nn.Module):
5     def __init__(self, linear):
6         super().__init__()
7         self.linear = linear
8         self.b = nn.Parameter(torch.zeros(self.
9             linear.out_features))
10
11     def forward(self, x):
12         ilora = self.b * x.sum(-1, keepdim=True)
13         return self.linear(x) + ilora

```

Listing 1. Implementation of 1LoRA in Python

### B. Ablations

#### B.1. Complementary methods

We investigate all possible combinations of BitFit, DiffFit and 1LoRA, including individual components (the biases  $\beta$ , the scaling factors  $\gamma$ , and the norms, see Table 1), in order to understand which modules impact fine-tuning most. We report experiments for MDE, fine-tuning DepthAnything from KITTI to NYU.

Interestingly, Figure 10 shows that, for the MDE case, only fine-tuning the normalization layers (norms, ●) leads to suboptimal results, whereas combining it with 1LoRA ● allows it to achieve similar results as 1LoRA ●, but faster. Fine-tuning the biases alone (i.e. BitFit, ●) achieves below state-of-the-art performance, but combining it with 1LoRA ● has the same effect as combining norms with 1LoRA ●. Other combinations seem to lead to no additional benefit, and only slow down the fine-tuning process. In general, any combination also increases the amount of trainable parameters, and 1LoRA already achieves the best results while using the least number of trainable parameters in this experiment.

#### B.2. Summation vs random compression

Here, we report ablation on the MDE (Table S1) task comparing 1LoRA with fixed summation and random compression. Results show that summation is overall slightly better in terms of performance or required time, suggesting that

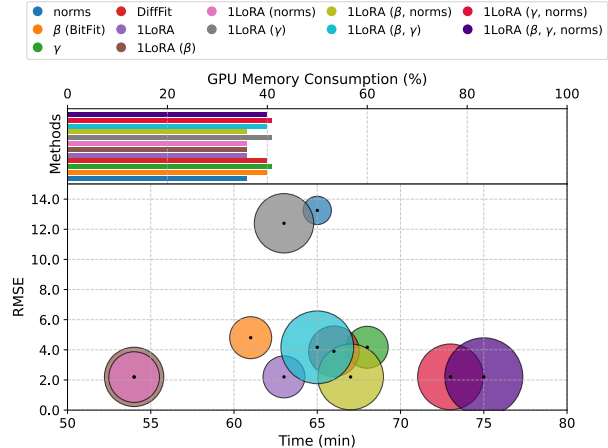


Figure 10. RMSE ( $\downarrow$ ) of DepthAnything model pre-trained on KITTI and fine-tuned to NYU, using all combinations of 1LoRA, BitFit and DiffFit, including individual components (biases  $\beta$ , scaling factors  $\gamma$  and normalization layers *norms*). Bubble size is proportional to the number of parameters. Bottom left (and smallest bubble) is better.

using a fixed summation compression is at least as good as using random compression.

Table S1. (MDE) 1LoRA: summation vs random compression.

dataset	compression	time [min] ( $\downarrow$ )	RMSE ( $\downarrow$ )
NYU $\rightarrow$ KITTI	summation	<b>63</b>	<b>2.203 <math>\pm</math> 0.0078</b>
	random	65	2.239 $\pm$ 0.0089
KITTI $\rightarrow$ NYU	summation	<b>65</b>	<b>0.238 <math>\pm</math> 0.0005</b>
	random	67	0.24 $\pm$ 0.0002
KITTI $\rightarrow$ DIODE Out	summation	<b>13</b>	<b>4.574 <math>\pm</math> 0.1217</b>
	random	22	4.573 $\pm$ 0.0582
NYU $\rightarrow$ DIODE Out	summation	42	<b>4.719 <math>\pm</math> 0.06</b>
	random	42	4.725 $\pm$ 0.06

#### B.3. Analysis of edge cases

To explore edge cases, we analyzed the error distribution and calibration of full fine-tuning, LoRA  $r = 1$  and 1LoRA for our classification experiments (Fig. 11). Contrary to expectation, full fine-tuning seems to be more heavy-tailed in terms of error distribution than the very low-rank PEFT methods. As for the calibration, again the very low-rank PEFT methods seem to slightly improve above full fine-

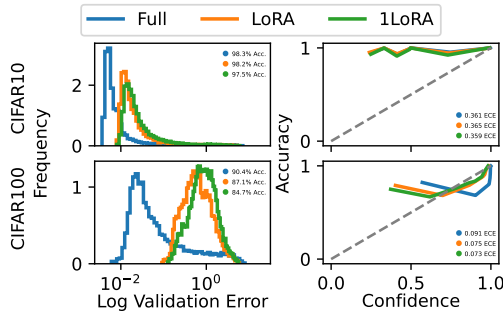


Figure 11. Error distribution and calibration for classification examples.

tuning, though the improvement – if any – is very much negligible.

### C. Mathematical reasoning visualization

We report visualization for LLaMA2 13b finetuned on 80% MetaMath, comparing 1LoRA and LoRA QKV. We use random samples from our generated validation split (20%) of MetaMath. We use black, olive and cyan boxes to highlight the input prompt (question and answer, where only the question is given as input and the answer serves as ground truth), the 1LoRA and LoRA QKV answers, respectively.

#### Question & Answer

### Instruction:

On Friday, a snowboard originally priced at \$100 was discounted 50%. On Monday, that sale price was reduced by 30%. In dollars, what is the price of the snowboard after the Monday reduction?

### Response:

The Friday sale price was  $\$100 \times 0.50 = \$50$ . The Monday reduction is 30% of \$50, so it is  $\$50 \times 0.30 = \$15$ . Therefore, the price of the snowboard after the Monday reduction is  $\$50 - \$15 = \$35$ . The answer is: 35.

#### 1LoRA

The Friday discount of 50% reduces the price to  $100 - 0.5(100) = 50$ . The Monday reduction of 30% reduces the price further to  $50 - 0.3(50) = 35$ . The answer is: 35

#### LoRA QKV

The snowboard was originally priced at \$100. The 50% discount brings the price down to  $\$100 \cdot 0.5 = \$50$ . The 30% reduction brings the price down further to  $\$50 \cdot 0.7 = \$35$ . Therefore, the price of the snowboard after the Monday reduction is  $\$35$ . The answer is: 35

#### Question & Answer

### Instruction:

George went to a movie theater to see a movie. He paid \$16 for the ticket and decided to buy some nachos. The nachos were half the price of the ticket. How much did George pay in total for the visit to the movie theater?

### Response:

George paid \$16 for the ticket. The nachos were half the price of the ticket, so they cost  $\$16/2 = \$8$ . In total, George paid  $\$16 + \$8 = \$24$  for the visit to the movie theater.

#### 1LoRA

The nachos were half the price of the ticket, so they cost  $\$16/2 = \$8$ . Therefore, George paid a total of  $\$16 + \$8 = \$24$  for the visit to the movie theater.

### 24

The answer is: 24

#### LoRA QKV

The nachos were half the price of the ticket, so they cost  $\$16/2 = \$8$  In total, George paid  $\$16 + \$8 = \$24$  for the visit to the movie theater.

### 24

The answer is: 24

#### Question & Answer

### Instruction:

If  $ab + bc + cd + da = 30$  and  $b + d = 5$ , what is the value of  $a + c$ ?

### Response:

We can rewrite the given equation as  $a(b+d) + c(b+d) = 30$ . Using the second equation, we substitute  $b + d = 5$  to get  $5a + 5c = 30$ . Dividing both sides by 5, we have  $a + c = 6$ . The answer is: 6

### 1LoRA

We can rewrite the equation  $ab + bc + cd + da = 30$  as  $(a + c)(b + d) = 30$ . Substituting  $b + d = 5$ , we have  $(a + c)(5) = 30$ . Dividing both sides by 5, we get  $a + c = \boxed{6}$ . The answer is: 6

### LoRA QKV

We can rewrite the given equation as  $a(b+d) + c(b+d) = 30$ . Using the given information that  $b + d = 5$ , we have  $2ab + 2cd = 30$ . We can rewrite this equation as  $ab + cd = 15$ . Since  $ab + bc + cd + da = 30$ , we can subtract  $ab + cd$  from both sides to get  $bc + da = 15$ . We can also rewrite this equation as  $bc + da = 15$ . Since  $b + d = 5$ , we can substitute this into the equation to get  $bc + da = 5$ . Now, we can solve for  $a + c$  by subtracting  $a$  and  $c$  from both sides:  $a + c = 5$ . Therefore, the value of  $a + c$  is  $\boxed{5}$ . The answer is: 5

### Question & Answer

### Instruction:

In 5 years, Raven will be 4 times as old as Phoebe. If Phoebe is currently  $x$  years old, Raven is 55 years old. What is the value of unknown variable  $x$ ?

### Response:

To solve this problem, we need to determine the value of  $x$ , which represents Phoebe's current age. Let's break down the information given: In 5 years, Raven will be 4 times as old as Phoebe. Raven is currently 55 years old. We can set up the equation as follows: Raven's age in 5 years =  $4 * \text{Phoebe's age in 5 years}$  (Raven's current age + 5) =  $4 * (\text{Phoebe's current age} + 5)$   $55 + 5 = 4 * (x + 5)$   $60 = 4 * (x + 5)$  Let's simplify and solve for  $x$ :  $60 = 4x + 20$  To isolate  $x$ , we subtract 20 from both sides of the equation:  $60 - 20 = 4x + 20 - 20$   $40 = 4x$  To solve for  $x$ , we divide both sides of the equation by 4:  $40 / 4 = 4x / 4$   $10 = x$  The value of  $x$  is 10.

#### 10

The answer is: 10.

### 1LoRA

To solve this problem, we need to determine the value of  $x$ , which represents Phoebe's current age. We know that in 5 years, Raven will be 4 times as old as Phoebe. Therefore, Raven's age in 5 years will be 4 times Phoebe's age in 5 years. Let's break down the information given: Raven's current age: 55 years Phoebe's current age:  $x$  years In 5 years, Raven's age:  $55 + 5 = 60$  years In 5 years, Phoebe's age:  $x + 5$  years We can set up the equation as follows: Raven's age in 5 years =  $4 * \text{Phoebe's age in 5 years}$   $60 = 4 * (x + 5)$  Let's simplify and solve for  $x$ :  $60 = 4x + 20$  To isolate  $x$ , we subtract 20 from both sides of the equation:  $60 - 20 = 4x + 20 - 20$   $40 = 4x$  To solve for  $x$ , we divide both sides of the equation by 4:  $40 / 4 = 4x / 4$   $10 = x$  The value of  $x$  is 10.

#### 10

The answer is: 10

### LoRA QKV

To solve this problem, we need to determine the value of  $x$ , which represents Phoebe's current age. We know that in 5 years, Raven will be 4 times as old as Phoebe. Therefore, Raven will be  $4x$  years old in 5 years. We also know that Raven is currently 55 years old. We can set up the equation as follows: Raven's age in 5 years =  $4 * \text{Phoebe's current age}$   $55 = 4 * x$  Let's simplify and solve for  $x$ :  $55 = 4x$  To isolate  $x$ , we divide both sides of the equation by 4:  $55 / 4 = 4x / 4$   $14 = x$  The value of  $x$  is 14.

#### 14

The answer is: 14

## D. Detailed results

We report tables (from S2 to S12) illustrating the detailed results shown in the main paper as bubble and bar plots. For each column reporting a result, we highlight the best value in **bold**, the second best in underlined and the third in *italic*.

method	time [min] ( $\downarrow$ )	RMSE ( $\downarrow$ )	memory [%] ( $\downarrow$ )	params ( $\downarrow$ )
MoRA <sub>1</sub>	69	<i>0.248 ± 0.00021</i>	<b>36</b>	<u>220.06K</u>
MoRA <sub>6</sub>	73	0.277 ± 0.00066	<b>36</b>	<b>217.44K</b>
BitFit	<b>64</b>	0.349 ± 0.00073	<u>40</u>	<i>221.18K</i>
DiffFit	78	0.314 ± 0.00078	<u>40</u>	321.54K
LoRA	71	0.25 ± 0.00087	<i>41</i>	393.22K
VeRA	69	<u>0.244 ± 0.00059</u>	<u>40</u>	221.28K
DoRA	104	0.249 ± 0.00088	54	614.40K
1LoRA	<u>65</u>	<b>0.238 ± 0.0005</b>	<b>36</b>	<i>221.18K</i>
All	83	0.211 ± 0.00064	56	335.32M

Table S2. Table for Figure 2a

method	time [min] ( $\downarrow$ )	RMSE ( $\downarrow$ )	memory [%] ( $\downarrow$ )	params ( $\downarrow$ )
MoRA <sub>1</sub>	67	2.509 ± 0.00455	<b>36</b>	<u>220.06K</u>
MoRA <sub>6</sub>	71	2.97 ± 0.00328	<b>36</b>	<b>217.44K</b>
BitFit	<b>61</b>	4.806 ± 0.00392	<u>40</u>	<i>221.18K</i>
DiffFit	66	3.904 ± 0.00539	<u>40</u>	321.54K
LoRA	66	<u>2.253 ± 0.00315</u>	<i>41</i>	393.22K
VeRA	68	<i>2.304 ± 0.00524</i>	<u>40</u>	221.28K
DoRA	100	2.346 ± 0.00275	54	614.40K
1LoRA	<u>63</u>	<b>2.203 ± 0.0078</b>	<b>36</b>	<i>221.18K</i>
All	85	1.916 ± 0.00303	56	335.32M

Table S3. Table for Figure 2b

method	time [min] ( $\downarrow$ )	RMSE ( $\downarrow$ )	memory [%] ( $\downarrow$ )	params ( $\downarrow$ )
MoRA <sub>1</sub>	52	4.704 ± 0.01546	<b>36</b>	<u>220.06K</u>
MoRA <sub>6</sub>	<u>14</u>	<i>4.637 ± 0.00449</i>	<b>36</b>	<b>217.44K</b>
BitFit	47	5.158 ± 0.00868	<u>40</u>	<i>221.18K</i>
DiffFit	50	4.955 ± 0.0064	<u>40</u>	321.54K
LoRA	51	<u>4.612 ± 0.00368</u>	<i>41</i>	393.22K
VeRA	24	4.797 ± 0.0181	<u>40</u>	221.28K
DoRA	74	4.644 ± 0.00405	54	614.40K
1LoRA	<b>13</b>	<b>4.574 ± 0.1217</b>	<b>36</b>	<i>221.18K</i>
All	55	4.59 ± 0.0784	56	335.32M

Table S4. Table for Figure 2c

method	time [min] ( $\downarrow$ )	RMSE ( $\downarrow$ )	memory [%] ( $\downarrow$ )	params ( $\downarrow$ )
MoRA <sub>1</sub>	51	5.067 $\pm$ 0.02273	<b>36</b>	220.06K
MoRA <sub>6</sub>	54	6.242 $\pm$ 0.0141	<b>36</b>	<b>217.44K</b>
BitFit	<u>47</u>	7.698 $\pm$ 0.00293	<u>40</u>	221.18K
DiffFit	<u>50</u>	7.017 $\pm$ 0.00998	<u>40</u>	321.54K
LoRA	67	4.94 $\pm$ 0.00776	<u>41</u>	393.22K
VeRA	51	<u>4.839 <math>\pm</math> 0.01025</u>	40	221.28K
DoRA	75	<u>4.918 <math>\pm</math> 0.00949</u>	54	614.40K
1LoRA	<b>42</b>	<b>4.719 <math>\pm</math> 0.066</b>	<b>36</b>	221.18K
All	61	4.53 $\pm$ 0.0278	56	335.32M

Table S5. Table for Figure 2d

method	time [min] ( $\downarrow$ )	Loss ( $\downarrow$ )	memory [%] ( $\downarrow$ )	params ( $\downarrow$ )
MoRA <sub>1</sub>	187	0.18 $\pm$ 0.000724	80	1.36M
MoRA <sub>6</sub>	198	0.19 $\pm$ 6.6e-05	91	<u>1.35M</u>
BitFit	<b>130</b>	0.22 $\pm$ 0.0047	<b>75</b>	1.36M
DiffFit	<u>132</u>	0.21 $\pm$ 0.00411	77	1.66M
LoRA	154	<b>0.15 <math>\pm</math> 0.00202</b>	93	2.50M
VeRA	159	<u>0.176 <math>\pm</math> 0.00192</u>	92	1.36M
1LoRA	152	<u>0.17 <math>\pm</math> 0.003768</u>	<u>76</u>	1.36M
LoRA (QKV)	<u>134</u>	0.2 $\pm$ 0.002273	79	<b>786.43K</b>

Table S6. Table for Figure 3

method	time [min] ( $\downarrow$ )	Loss ( $\downarrow$ )	memory [%] ( $\downarrow$ )	params ( $\downarrow$ )
BitFit	<u>66</u>	0.19 $\pm$ 0.001738	<b>91</b>	<u>2.13M</u>
DiffFit	<u>71</u>	<u>0.18 <math>\pm</math> 0.003334</u>	93	2.58M
1LoRA	77	<b>0.14 <math>\pm</math> 0.00152</b>	<u>92</u>	<u>2.13M</u>
LoRA (QKV)	<b>61</b>	<u>0.17 <math>\pm</math> 0.005107</u>	<u>92</u>	<b>1.23M</b>

Table S7. Table for Figure 4

method	time [h] ( $\downarrow$ )	FID ( $\downarrow$ )	memory [%] ( $\downarrow$ )	params ( $\downarrow$ )
MoRA <sub>1</sub>	<b>9.6</b>	13.9004	67	604.41K
MoRA <sub>6</sub>	23.5	14.0791	87	<b>596.52K</b>
BitFit	<u>13.0</u>	13.8074	<u>66</u>	<u>603.68K</u>
DiffFit	14.3	<u>13.6644</u>	87	1.09M
LoRA	19.5	13.7621	79	864.03K
VeRA	<u>12.0</u>	<u>13.1129</u>	87	<u>603.82K</u>
1LoRA	17.0	<b>12.8218</b>	<b>65</b>	<u>603.68K</u>

Table S8. Table for Figure 5

method	time [min] ( $\downarrow$ )	Accuracy ( $\uparrow$ )	memory [%] ( $\downarrow$ )	params ( $\downarrow$ )
MoRA <sub>1</sub>	30	96.55 $\pm$ 0.0465	<b>6</b>	333.36K
MoRA <sub>6</sub>	31	97.5 $\pm$ 0.0544	<b>6</b>	<b>328.13K</b>
BitFit	<b>16</b>	97.0 $\pm$ 0.055	<b>6</b>	<u>331.78K</u>
DiffFit	18	97.9 $\pm$ 0.0585	8	817.15K
LoRA	33	<u>98.1 <math>\pm</math> 0.058</u>	<u>7</u>	663.55K
VeRA	29	97.7 $\pm$ 0.1631	8	332.06K
DoRA	51	<b>98.2 <math>\pm</math> 0.0866</b>	11	995.33K
1LoRA	<u>17</u>	97.7 $\pm$ 0.186	<b>6</b>	<u>331.78K</u>
1LoRA (norms)	18	98.0 $\pm$ 0.1397	<b>6</b>	485.38K
All	22	98.3 $\pm$ 0.0755	10	343.19M

Table S9. Table for Figure 6

method	time [min] ( $\downarrow$ )	Accuracy ( $\uparrow$ )	memory [%] ( $\downarrow$ )	params ( $\downarrow$ )
MoRA <sub>1</sub>	31	82.7 $\pm$ 0.3253	<b>6</b>	333.36K
MoRA <sub>6</sub>	30	81.2 $\pm$ 0.1875	<b>6</b>	<b>328.13K</b>
BitFit	<b>16</b>	83.0 $\pm$ 0.204	<b>6</b>	<u>331.78K</u>
DiffFit	18	87.9 $\pm$ 0.0775	8	817.15K
LoRA	34	<b>88.4 <math>\pm</math> 0.2173</b>	<u>7</u>	663.55K
VeRA	29	85.3 $\pm$ 0.1723	8	332.06K
DoRA	50	87.3 $\pm$ 0.2245	11	995.33K
ILoRA	<u>17</u>	85.0 $\pm$ 0.567	<b>6</b>	<u>331.78K</u>
ILoRA (norms)	18	<u>88.2 <math>\pm</math> 0.314</u>	<b>6</b>	485.38K
All	22	90.9 $\pm$ 0.3251	10	343.19M

Table S10. Table for Figure 7

method	time [min] ( $\downarrow$ )	RMSE ( $\downarrow$ )	memory [%] ( $\downarrow$ )	params ( $\downarrow$ )
norms	65	13.25 $\pm$ 0.004751	<b>36</b>	<b>100.35K</b>
$\beta$ (BitFit)	<u>61</u>	4.806 $\pm$ 0.003922	<u>40</u>	<u>221.18K</u>
$\gamma$	68	4.17 $\pm$ 0.00482	<u>41</u>	<u>221.18K</u>
DiffFit	66	3.904 $\pm$ 0.00539	<u>40</u>	321.54K
1LoRA	63	2.203 $\pm$ 0.0078	<b>36</b>	<u>221.18K</u>
1LoRA ( $\beta$ )	<b>54</b>	2.203 $\pm$ 0.005498	<b>36</b>	442.37K
1LoRA (norms)	<b>54</b>	<b>2.19 <math>\pm</math> 0.00655</b>	<b>36</b>	321.54K
1LoRA ( $\gamma$ )	63	12.399 $\pm$ 0.006093	<u>41</u>	442.37K
1LoRA ( $\beta$ , norms)	67	<u>2.2 <math>\pm</math> 0.0051438</u>	<b>36</b>	542.72K
1LoRA ( $\beta$ , $\gamma$ )	65	4.166 $\pm$ 0.008069	<u>40</u>	663.55K
1LoRA ( $\gamma$ , norms)	73	2.206 $\pm$ 0.001668	<u>41</u>	542.72K
1LoRA ( $\beta$ , $\gamma$ , norms)	75	2.206 $\pm$ 0.002445	<u>40</u>	763.90K

Table S11. Table for Figure 10



method	time [min] ( $\downarrow$ )	Accuracy ( $\uparrow$ )	memory [%] ( $\downarrow$ )	params ( $\downarrow$ )
norms	<b>16</b>	$97.7 \pm 0.0395$	<b>6</b>	<b>153.60K</b>
$\beta$ (BitFit)	<b>16</b>	$97.0 \pm 0.055$	<b>6</b>	<u>331.78K</u>
$\gamma$	<u>17</u>	$97.3 \pm 0.0457$	<u>8</u>	<u>331.78K</u>
DiffFit	<i>18</i>	<u><math>97.9 \pm 0.0585</math></u>	<u>8</u>	817.15K
1LoRA	<u>17</u>	$97.7 \pm 0.186$	<b>6</b>	<u>331.78K</u>
1LoRA ( $\beta$ )	<i>18</i>	$97.8 \pm 0.1245$	<b>6</b>	663.55K
1LoRA (norms)	<i>18</i>	<b><math>98.0 \pm 0.139</math></b>	<b>6</b>	485.38K
1LoRA ( $\gamma$ )	19	$97.8 \pm 0.1043$	<u>8</u>	663.55K
1LoRA ( $\beta$ , norms)	19	<u><math>97.9 \pm 0.078</math></u>	<b>6</b>	817.15K
1LoRA ( $\beta$ , $\gamma$ )	21	$97.8 \pm 0.1297$	<b>6</b>	995.33K
1LoRA ( $\gamma$ , norms)	20	<b><math>98.0 \pm 0.0837</math></b>	<u>8</u>	817.15K
1LoRA ( $\beta$ , $\gamma$ , norms)	22	<b><math>98.0 \pm 0.109</math></b>	<u>8</u>	1.15M

Table S12. Table for Figure 8