

Outlier-weighted Layerwise Sampling for LLM Fine-tuning

Pengxiang Li^{1,*} Lu Yin^{2,3,*} Xiaowei Gao⁴ Shiwei Liu^{5,3,†}

¹Dalian University of Technology ²University of Surrey ³Eindhoven University of Technology

⁴University College London ⁵University of Oxford

shiwei.liu@maths.ox.ac.uk*

Abstract

The rapid advancements in Large Language Models (LLMs) have revolutionized various natural language processing tasks. However, the substantial size of LLMs presents significant challenges in training or fine-tuning. While parameter-efficient approaches such as low-rank adaptation (LoRA) have gained popularity, they often compromise performance compared to full-rank fine-tuning. In this paper, we propose *Outlier-weighted Layerwise Sampling (OWS)*, a new memory-efficient fine-tuning approach, inspired by the layerwise outlier distribution of LLMs. Unlike LoRA, which adds extra adapters to all layers, OWS strategically assigns higher sampling probabilities to layers with more outliers, selectively sampling only a few layers and fine-tuning their pre-trained weights. To further increase the number of fine-tuned layers without a proportional rise in memory costs, we incorporate gradient low-rank projection, further boosting the approach’s performance. Our extensive experiments across various architectures, including LLaMa2, and Mistral, demonstrate that OWS consistently outperforms baseline approaches, including full fine-tuning. Specifically, it achieves up to a 1.1% average accuracy gain on the Commonsense Reasoning benchmark, a 3.0% improvement on MMLU, and a notable 10% boost on MT-Bench, while being more memory efficient. OWS allows us to fine-tune 7B LLMs with only 21GB of memory. Our code is available at <https://github.com/pixeli99/OWS>.

1 Introduction

The rapid advancements in AI driven by Large Language Models (LLMs) have fundamentally transformed how people work and communicate. The impressive language capabilities of LLMs enable a single model to handle various tasks simultaneously, including but not limited to natural language

understanding (Brown et al., 2020; Touvron et al., 2023), text generation (Kocoń et al., 2023; Anil et al., 2023), machine translation (Jiao et al., 2023), and programming (Surameery and Shakor, 2023; Tian et al., 2023). However, the massive size of LLMs presents significant challenges for practical applications and deployment.

To address these challenges, various parameter-efficient fine-tuning (PEFT) approaches have been proposed, including prompt tuning (Lester et al., 2021; Liu et al., 2021a), adapters (Houlsby et al., 2019; He et al., 2021), and low-rank adaptation (LoRA) (Hu et al., 2021; Dettmers et al., 2024). These approaches enable the fine-tuning of pre-trained LLMs with substantially fewer trainable parameters, making LLM fine-tuning more feasible in practice. Among these, LoRA (Hu et al., 2021) stands out for its re-parameterization technique of the pre-trained weight matrix $W \in \mathbb{R}^{m \times n}$, expressed as $W_0 + AB$, where $A \in \mathbb{R}^{m \times r}$, $B \in \mathbb{R}^{r \times n}$, and $r \ll \min(m, n)$. By fine-tuning only the low-rank adaptor AB while keeping the pre-trained weight W_0 frozen, LoRA significantly reduces the memory usage and computational costs associated with fine-tuning LLMs, rapidly becoming the preferred method for such tasks. Despite its efficiency, recent research has highlighted the inferior performance of low-rank reparameterization compared to full-rank updates in both fine-tuning scenarios (Xia et al., 2024; Biderman et al., 2024) and pre-training contexts (Lialin et al., 2023b; Zhao et al., 2024). These findings underscore the need for further exploration into balancing training efficiency with model performance, particularly in the context of large-scale language models.

Recently, Layerwise Importance Sampled AdamW (LISA) (Pan et al., 2024) has emerged as a promising alternative for LLM fine-tuning, integrating the concept of importance sampling (Kloek and Van Dijk, 1978; Zhao and Zhang, 2015) into the fine-tuning process. Unlike methods that add

*Equal contribution. †Corresponding author.

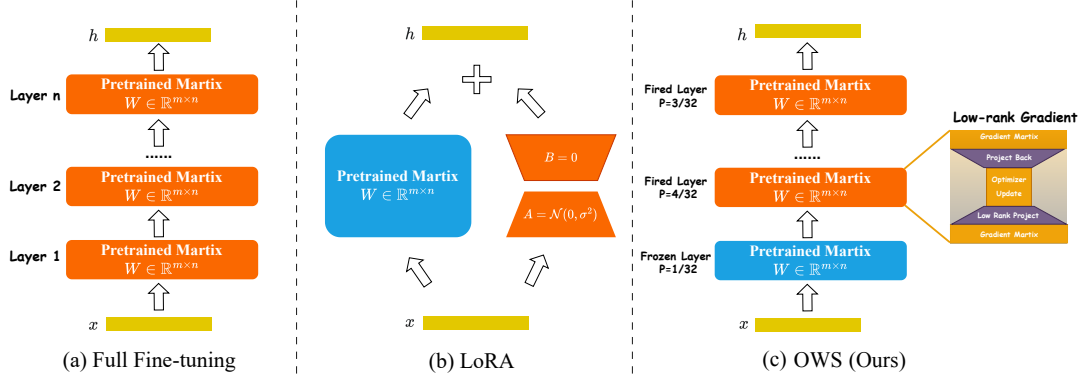


Figure 1: The comparison among Full Fine-tuning, training with LoRA, and OWS. Blue modules are frozen, while orange modules are activated. OWS non-uniformly samples layers to fine-tune models with low-rank gradients.

adapters to all layers, LISA selectively samples a small subset of layers and directly fine-tunes their pre-trained weights, demonstrating notable performance improvements over LoRA. However, our investigation reveals two limitations of LISA:

- ★ LISA employs random sampling of layers for fine-tuning, which results in suboptimal performance due to the varying importance of layers in LLMs. To illustrate this, we demonstrate that random sampling underperforms compared to a simple baseline—monotonic decreasing sampling probabilities from top to bottom layers—as shown in Table 1.
- ★ The fine-tuning of sampled layers is conducted in a full-rank manner, leading to substantial memory overhead as the number of sampled layers increases. While fine-tuning accuracy improves with the inclusion of more sampled layers, this improvement comes at the cost of escalating memory usage, as detailed in Table 8.

Overview. To address these limitations, we introduce Outlier-weighted Layerwise Sampling (OWS), a novel layerwise sampling approach to fine-tune LLMs. OWS leverages the unique characteristic of LLMs where certain features and weights—referred to as outliers—have significantly larger magnitudes than the rest (Kovaleva et al., 2021; Puccetti et al., 2022; Dettmers et al., 2022; Yin et al., 2024; Lu et al., 2024). Our rationale is that layers with more outliers play a more active role in learning the data distribution of the text corpus, as they have received larger gradients during pre-training. Therefore, we assign higher sampling probabilities to layers with a greater concentration of outliers, in a way that the features in these layers will be leveraged more frequently to

adapt to downstream datasets. To further reduce the memory costs of fine-tuning, we update the optimization status of sampled layers in a low-rank subspace, which allows us to increase the number of fine-tuned layers without a proportional rise in memory costs.

Built upon these techniques, OWS achieves a large performance boost to LISA. Our extensive experiments across commonly used architectures including LLaMa2 (Touvron et al., 2023) and Mistral (Jiang et al., 2023) demonstrate that OWS consistently outperforms its baseline approaches including full-parameter fine-tuning. OWS achieves up to a 1.1% average accuracy gain on the Commonsense Reasoning benchmark, a 3.0% improvement on MMLU, and a notable 10% boost on MT-Bench.

2 Background and Motivation

In this section, we first introduce LISA’s algorithm and then present our findings of two key limitations of LISA: the shortcomings of its sampling approach and the significant memory overhead associated with the sampled layers.

2.1 Background: LISA

Pan et al. (2024) conducted an in-depth analysis of LoRA’s training dynamics across layers and revealed an unusual skew in the distribution of layerwise weight norms, particularly towards the top layer and/or the bottom layer¹, where the norms are significantly larger compared to other layers. Building upon this insight, the authors proposed LISA, a novel fine-tuning approach for LLMs, which incorporates the concept of importance sampling (Kloek and Van Dijk, 1978; Zhao and Zhang,

¹Please note that in LISA, the terms ‘top’ and ‘bottom’ layers refer to the embedding layer and the LLM head layer, respectively, rather than the first and last Transformer blocks.

2015) into the fine-tuning process. In LISA, layers of the base model are sampled to be unfrozen during training based on a prescribed probability, with the exception of the top and bottom layers, which remain activated throughout the process. Given a network with N_L layers, the sampling probability of layer ℓ is given as follows:

$$p_\ell = \begin{cases} 1.0, & \text{if } \ell = 1 \text{ or } \ell = N_L, \\ \gamma/N_L & \text{else.} \end{cases} \quad (1)$$

where γ controls the expected number of unfrozen layers during optimization. Since LISA does not require additional adaptors and only fine-tunes an expected γ layers, it notably reduces the memory usage of LLM fine-tuning.

2.2 Limitations of LISA

While demonstrating promising results, we observe that the LISA algorithm inherently has two shortcomings that constrain its memory-performance trade-off:

i. The middle layers of LISA are sampled uniformly, which can result in suboptimal performance. To verify this point, we conduct a small experiment where we replace the uniform sampling with a very simple baseline, i.e. monotonic decreasing sampling, where the sample probability is monotonically decreasing from shallow layers to deep layers (noted as LISA-D). Table 1 shows that this simple sampling method often outperforms uniform sampling, verifying our concern.

ii. The sampled layers of LISA are fine-tuned in a full-rank manner, causing a significant memory increase as the number of sampled layers increases. To illustrate this, we fine-tune LLaMA2-7B on the GSM8K training set and report the GSM8K score and memory usage of LISA with various numbers of sampled layers γ , as shown in Table 8. The memory requirement of LISA rises significantly from 23GB to 36GB as γ increases from 1 to 12. Similarly, the performance improves consistently with the increase in sampled layers. Since sampling more layers results in stronger fine-tuning performance, it is crucial to reduce the associated memory overhead as the number of sampled layers grows.

3 Outlier-weighted Layerwise Sampling (OWS)

In this section, we introduce our approach, Outlier-weighted Layerwise Low-Rank Projection (OWS).

We will discuss the underlying rationales, present preliminary results, and detail the algorithm design.

The above findings shed light on a principle for designing non-uniform layerwise sampling for LLM fine-tuning: layers with higher outlier ratios should be prioritized during the fine-tuning process. This forms the foundation of our proposed method, Outlier-weighted Layerwise Low-Rank Projection (OWS), which we will present in detail.

Outlier-Weighted Sampling (OWS). Although LISA-D achieves good performance, it is more desirable to seek a more principled approach to determine the layerwise sampling probability. In the context of LLMs, we get inspiration from the unique characteristic of LLMs, outliers, defined as features and weights exhibiting significantly larger magnitudes compared to the majority of others (Kovaleva et al., 2021; Puccetti et al., 2022; Dettmers et al., 2022). It has been widely demonstrated that removing outliers significantly degrades the capacity of LLMs (Dettmers et al., 2022).

Our motivation stems from the crucial role outliers play in preserving LLM performance (Yin et al., 2024). We hypothesize that layers with more outliers likely contain more essential information, as they have received larger gradients during training. Therefore, we assign higher sampling probabilities to layers with more outliers during fine-tuning, leading to a substantial improvement in performance. To formulate, let us consider the input of a layer as \mathbf{X} with dimensions $(N \times L, C_{in})$, where N and L represent the batch and sequence dimensions, respectively; and the weight matrix \mathbf{W} has dimensions (C_{out}, C_{in}) . Outlier score of weight \mathbf{W}_{ij} is computed as $\mathbf{A}_{ij} = \|\mathbf{X}_j\|_2 \cdot |\mathbf{W}_{ij}|$. Here, $\|\mathbf{X}_j\|_2$ is the ℓ_2 norm of input feature connected to \mathbf{W}_{ij} .

We first calculate the layerwise outlier distribution of a N_L -layer as $[D_1, D_2, \dots, D_{N_L}]$, where D_ℓ characterizes the outlier ratio of layer ℓ :

$$D_\ell = \frac{\sum_{i=1}^{C_{out}} \sum_{j=1}^{C_{in}} \mathbb{I}(\mathbf{A}_{ij}^\ell > \tau \cdot \bar{\mathbf{A}}^\ell)}{C_{in} C_{out}}, \quad (2)$$

where $\bar{\mathbf{A}}^\ell$ is the mean of \mathbf{A}^ℓ and $\mathbb{I}(\cdot)$ is the indicator function, returning 1 if \mathbf{A}_{ij}^ℓ is larger than $\tau \cdot \bar{\mathbf{A}}^\ell$, else 0. The layerwise outlier distribution essentially counts up weights whose outlier score is τ^2 times greater than that layer’s average outlier score. Larger D means more outliers are presented in the corresponding layer. The sampling probability p_ℓ

²We empirically find $\tau = 13$ consistently works well and choose it for all experiments in this paper.

Table 1: Fine-tuning performance of LLaMA2-7B with various dataset. The results are averaged under three random seeds.

Model	Method	BoolQ	PIQA	SIQA	HellaSwag	WinoGrande	OBQA	Average
Llama2-7B	LISA	82.0	79.9	33.5	59.7	79.6	38.8	62.25
Llama2-7B	LISA-D	85.1	79.9	33.8	59.8	79.7	38.4	62.78

of layer ℓ is then calculated as $p_\ell = \gamma D_\ell / \sum_{i=1}^{N_L} D_i$, where γ is the hyperparameter inherited from LISA to control the expected number of unfreeze layers during optimization. At each iteration, only the sampled layers will be fine-tuned, while the remaining layers are kept frozen. The visualization of layerwise outlier distribution of OWS is illustrated in Figure 2.

Can the weights in the sampled layer be further updated in a low-rank subspace? Outlier-weighted sampling addresses our first research question: how to optimally sample layers for sampling-based LLM fine-tuning. To tackle the second issue of the substantial memory cost associated with an increasing number of unfrozen layers, we propose to integrate outlier-weighted sampling with gradient low-rank training. In this approach, the sampled layers are updated in a low-rank manner (Zhao et al., 2024). Specifically, for each sampled layer, the gradient matrix is projected into a low-rank subspace using Singular Value Decomposition (SVD). The optimizer states are subsequently updated in the corresponding low-rank subspace with a rank level of r , significantly reducing the memory cost of optimization. We update the gradient subspace every 200 iterations to better capture the dynamic trajectory of fine-tuning. The above two innovations significantly boost the memory efficiency of OWS, unlocking the performance-memory trade-off of sampling-based fine-tuning. At the macro level, we dynamically sample a limited number of layers to fine-tune at each iteration. At the micro level, each sampled layers are updated with low-rank gradients.

Since the sampled layers are updated in the low-rank subspace, we can efficiently increase the number of sampled layers γ with only a marginal increase in memory cost compared to LISA. Additionally, as we sample only a few layers at each fine-tuning iteration, we can increase the rank levels r without significantly raising the memory requirements compared to LoRA. Memory usage analysis is given in Section 4.3. We perform a small search and find that $\gamma = 5$ and $r = 128$ consistently give us robust performance across models and downstream tasks. Therefore, we choose $\gamma = 5$ and

$r = 128$ as our default settings. We present our algorithm in Algorithm 1.

4 Experiments

In this section, we conduct extensive experiments³ to evaluate the effectiveness of OWS on multiple fine-tuning tasks. Details are provided below.

4.1 Experimental Setup

We choose multiple open-source LLMs that are widely used in research and practice, such as LLaMa2-7B (Touvron et al., 2023) and Mistral-7B (Jiang et al., 2023).

Fine-tuning Tasks. We choose an extensive range of fine-tuning tasks aiming to provide a thorough evaluation of OWS. Our fine-tuning tasks cover three categories: (i) **Commonsense Reasoning** (Hu et al., 2023), which includes 8 reasoning tasks including. (ii) **MT-Bench** (Zheng et al., 2024), a challenging multi-turn question set to assess the conversational and instruction-following abilities of models. We apply GPT-3.5-turbo and GPT-4o as the judge for MT-Bench; (iii) **MMLU** (Hendrycks et al., 2020), a massive multitask test consisting of multiple-choice questions from various branches of knowledge. We adopt the 5-shot setting for MMLU. For Commonsense Reasoning, all models are first fine-tuned on commonsense170k and then evaluated separately on different tasks, following Hu et al. (2023); For MT-Bench, we first fine-tune models on the Alpaca GPT-4 dataset (Peng et al., 2023) and then evaluate on MT-Bench following LISA. The results of MMLU are fine-tuned on the auxiliary training dataset and then evaluated on MMLU with 5 shots.

PEFT Baselines. We mainly consider four state-of-the-art baselines that are closely related to our approach: (i) **Full fine-tuning (Full FT)**: all parameters of pre-trained models are fine-tuned. Weights, gradients, and optimization states are maintained with full rank; (ii) **LoRA** (Hu et al., 2021): LoRA introduces additional low-rank adaptors and only fine-tunes adaptors, while maintaining pre-trained

³Our repository is built on top of LMFlow: <https://github.com/OptimalScale/LMFlow>

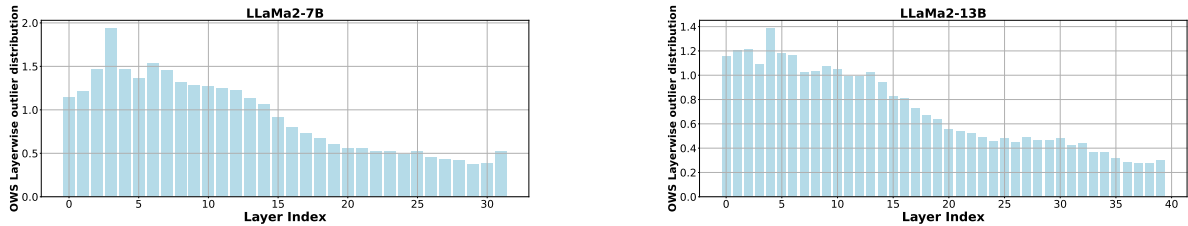


Figure 2: OWS Layerwise outlier distribution of LLaMa2 of Equation 2. The Y-axis is presented in percentage. Higher values mean higher outlier ratios.

Table 2: Fine-tuning performance of LLaMa2-7B and Mistral-7B with various approaches on commonsense reasoning datasets. The results are averaged under three random seeds.

Method	Mem.	BoolQ	PIQA	SIQA	HellaSwag	WinoGrande	ARC-e	ARC-c	OBQA	Avg.
LLaMa2-7B										
Full FT	61G	87.3	79.5	32.7	56.7	80.2	78.5	49.0	40.8	63.1
LoRA	26G	79.7	79.7	34.4	59.9	79.8	79.5	49.7	36.6	62.4
GaLore	36G	81.8	79.4	32.9	60.7	79.6	79.8	49.4	37.6	62.7
LISA	24G	82.0	79.9	33.5	59.7	79.6	80.4	51.1	38.8	63.1
OWS	23G	85.4	80.7	34.2	60.3	82.2	80.6	51.0	39.1	64.2
Mistral-7B										
Full FT	61G	86.5	84.3	33.5	65.1	87.1	83.8	57.5	41.2	67.4
LoRA	26G	87.2	81.0	33.7	62.9	83.3	82.2	54.2	37.0	65.2
GaLore	36G	84.8	82.5	34.4	63.5	85.6	82.5	53.9	37.8	65.6
LISA	24G	84.7	82.9	33.4	64.2	85.8	83.4	54.4	40.5	66.2
OWS	23G	87.8	83.9	34.0	66.4	85.6	84.1	57.9	40.4	67.5

weights frozen during training; (iii) **GaLore** (Zhao et al., 2024): pre-trained LLMs are fine-tuned with low-rank gradient projection. We follow (Zhao et al., 2024) and set the rank level to 8 for both GaLore and LoRA in all fine-tuning tasks; (iv) **LISA** (Pan et al., 2024): LISA is a sampling-based LLM fine-tuning method, which by default samples 2 layers to fine-tune with full rank at each iteration. GaLore and LISA directly fine-tune pre-trained weights without additional adaptors.

Hyperparameter Tuning. Regarding the hyperparameters of the baselines, we have conducted extensive hyperparameter tuning for all baselines with LLaMa2-7B and reported the results with the best ones. For Mistral-7B, we directly use the best hyperparameters of LLaMa2-7B. Specifically, for the learning rate, we performed a hyperparameter sweep over $[1e-4, 3e-4, 7e-5, 5e-5, 1e-5, 5e-6]$ for each method. For GaLore, we tested several update frequencies for the subspace $[50, 100, 200, 500]$ and found that 200 works best, consistent with GaLore’s reports. To ensure a fair comparison, we followed GaLore’s approach and set the rank level to 8 for GaLore and LoRA, resulting in approximately 24GB of memory usage for all methods. Additionally, we thoroughly analyzed the effect of two hyperparameters, such as rank level and sampled layers, as shown in Figure 3, where

our approach consistently demonstrates superior memory benefits. More configurations details are reported in Appendix D.

4.2 Experimental Results

In this section, we present the empirical results of OWS in comparison to other baseline methods.

Commonsense Reasoning Benchmark. We first evaluate with 8 commonsense reasoning tasks. The results are reported in Table 2. Overall, OWS consistently outperforms Full FT and other PEFT baselines by a large margin across various LLMs, demonstrating the superiority of OWS in LLM fine-tuning. We summarize our key observations below:

① **OWS approaches significantly outperform other efficient fine-tuning approaches by a large margin.** OWS consistently outperforms its layerwise sampling baseline, LISA, on nearly all tasks with LLaMa2-7B, delivering an average of 1.1% performance gain.

② **OWS outperforms full fine-tuning across tasks on LLaMa.** We can observe that OWS can achieve better performance than full fine-tuning with all models. LISA can match the performance of full fine-tuning for LLaMa models, whereas GaLore and LoRA perform no better than full fine-tuning. However, only OWS is able to match the performance of full fine-tuning with Mistral-7B and all other baselines fail to do so. This result

Table 3: Fine-tuning performance of LLaMa2-7B with various approaches on MT-Bench using GPT-3.5-turbo as a judge. The results are averaged under three random seeds.

Method	Writing	Roleplay	Reasoning	Math	Coding	Extraction	STEM	Humanities	Avg.
Full-FT	7.11	8.11	4.90	2.85	3.75	6.50	7.80	8.10	6.14
LoRA	7.21	7.05	4.95	3.25	3.90	5.70	7.90	7.65	5.95
GaLore	7.05	7.79	3.55	2.89	3.15	6.25	8.30	7.63	5.83
LISA	6.75	7.35	4.35	3.00	3.85	6.85	7.74	7.47	5.92
OWS	8.00	7.65	4.95	3.25	4.15	7.45	8.25	8.45	6.52

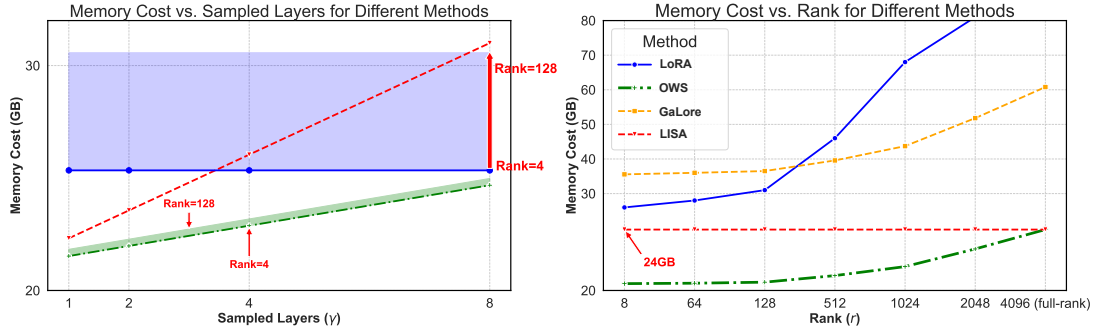


Figure 3: Fine-tuning memory usage of using various with LLaMa2-7B. **Left:** varying sampled layers. In this scenario, we also vary the rank of LoRA and OWS from 4 to 128 to provide a comprehensive analysis. OWS consistently demonstrates superior memory efficiency across all configurations. Notably, LISA’s memory advantage over LoRA diminishes as the number of sampled layers increases. **Right:** varying ranks. The sampled layer of LISA and OWS is set as $\gamma = 2$.

suggests that LLMs gain greater benefits by leveraging features within important layers rather than uniformly distributing resources across all layers for fine-tuning.

MT-Bench. We next evaluate OWS on a more comprehensive benchmark, MT-Bench, featuring 80 high-quality, multi-turn questions designed to assess LLMs on 8 common categories. Results are presented in Table 3. We can observe that the benefits of OWS over other PEFT approaches are more pronounced. Using GPT-3.5-turbo as a judge, all other baselines fail to match the performance of full fine-tuning on MT-Bench with scores below 6.0, whereas OWS outperforms the full fine-tuning by a large margin. To be specific, OWS significantly boosts the average score of LISA from 5.92 to 6.52.

Table 4: Mean score of LLaMA-2-7B on MT-Bench over three seeds. The results are averaged under three random seeds.

Judge	Full-FT	LoRA	GaLore	LISA	OWS
GPT-3.5-turbo	6.14	5.95	5.83	5.92	6.52
GPT-4o	4.91	4.58	4.73	4.81	5.10

The performance trend when using GPT-4 is very similar to that of GPT-3.5-turbo, although the scores evaluated by GPT-4 are generally lower. Notably, only OWS outperforms full fine-tuning, achieving a higher score over full fine-tuning.

MMLU Benchmark. To draw a more solid con-

clusion, we also test another widely used benchmark, i.e., MMLU. The results are shown in Table 5. Our findings highlight that OWS consistently outperforms Full FT, while other PEFT methods fall short of dense fine-tuning. Specifically, OWS achieves an average score of 52.6, demonstrating significant improvements across various domains such as Humanities, STEM, Social Sciences, and Others. These results underscore OWS’s efficacy beyond full fine-tuning while maintaining superior memory efficiency.

Table 5: Fine-tuning performance of LLaMa2-7B with various approaches on MMLU benchmark. The results are averaged under three random seeds.

Method	Humanities	STEM	Social.	Other	Avg.
Full-FT	49.9	41.7	57.5	57.0	51.5
LoRA	46.1	40.8	56.6	56.2	49.9
GaLore	45.4	41.7	55.8	56.0	49.7
LISA	44.9	41.2	54.7	57.6	49.6
OWS	49.8	42.1	58.6	59.7	52.6

GSM8K. We extend our evaluation to compare OWS with two recent memory-efficient fine-tuning methods, **HiFT** (Liu et al., 2024b) and **MeZO** (Malladi et al., 2023), on the GSM8K benchmark. Table 6 indicates that OWS achieves the best accuracy, surpassing HiFT by +1.6 percentage points (pp) and MeZO by +2.5 pp under the same model size.

Table 6: GSM8K accuracy (%) of LLaMA2-7B with different memory-efficient fine-tuning strategies.

Method	Model	GSM8K
HiFT (best: RAN)	LLaMA2-7B	20.3
MeZO	LLaMA2-7B	19.4
OWS ($r=128, \gamma=2$)	LLaMA2-7B	21.9

Generalisability to Newer Architectures. To examine scalability, we further fine-tune the recent Qwen2.5-7B model on GSM8K and compare OWS with DoRA (Liu et al., 2024a) and LISA. As shown in Table 7, OWS attains the highest score of 83.7%, outperforming DoRA by +2.5 pp and LISA by +4.0 pp, indicating strong generalisation to state-of-the-art LLMs.

Table 7: GSM8K accuracy (%) on Qwen2.5-7B.

Method	Model	GSM8K
LISA	Qwen2.5-7B	79.7
DoRA	Qwen2.5-7B	81.2
OWS	Qwen2.5-7B	83.7

4.3 Memory Efficiency of OWS

Thanks to its layerwise sampling and low-rank characteristics, OWS significantly improves the memory efficiency of LLM fine-tuning. To verify, we report the memory cost of various approaches when used to fine-tune LLaMa2-7B, with a token batch size of 1 in Figure 3.

On the one hand, the low-rank nature of OWS allows us to unfreeze more layers without a substantial increase in memory cost compared to LISA. As illustrated in Figure 3-Left, when increasing γ from 1 to 8, LISA exhibits a notable memory growth from 23GB to 32GB, whereas OWS’s memory cost slightly increases from 21GB to 25GB. Compared to LoRA with $r = 4$, OWS facilitates training with a much higher rank ($r = 128$) while still maintaining a lower memory cost. On the other hand, Figure 3-Right demonstrates that OWS enables high-rank training without significantly compromising memory efficiency, in stark contrast to LoRA. It is important to note that we do not utilize the layer-wise weight update technique used in GaLore for the memory measurement, hence the memory cost of GaLore is higher than reported in GaLore.

We further break down the memory usage during LLM fine-tuning, presenting the results in Figure

4-Left. For this analysis, the number of fine-tuned layers γ is set to 2 for both LISA and OWS, and rank level r is set to 8 for both LoRA and OWS. LoRA incurs a substantial activation memory cost, although its optimizer and gradient memory requirements are relatively small. In contrast, LISA’s optimizer memory cost is large because each layer is trained in full rank, yet it benefits from a small activation memory cost. OWS effectively combines the advantages of both methods, inheriting the small activation memory of LISA while significantly reducing the optimizer memory requirement.

4.4 Superiority of OWS under Varying Hyperparameters Over LISA

The primary hyperparameters of LISA, GaLore, and OWS are the number of fine-tuned layers γ , and the rank level within each layer r . To evaluate their effect on the performance of different approaches, we vary these two hyperparameters and report the results in Table 8. We set $\gamma = 32$ for GaLore and $r = \text{‘full rank’}$ for LISA as their default. We see that GaLore’s performance does improve as rank levels, having the lowest score across most cases. Notably, OWS significantly reduces the memory cost compared to LISA alone—reducing from 36G to 27G with $r = \text{full}, \gamma = 12$ —while achieving a significant improvement of 6.1.

4.5 OWS Serves as A Better Layerwise Important Metric than Others

OWS serves as a better layer-wise importance metric than previous ones. We compare OWS with other layerwise importance scores for sampling-based fine-tuning, including Uniform (Pan et al., 2024), Relative Magnitude (RM) (Samragh et al., 2023) and Block Influence (BI) (Men et al., 2024) in Table 9. OWS consistently performs better than other layer importance scores. Note that reversing OWS gives us the worse performance as shown in Appendix A.

4.6 Memory Usage Breakdown and Training Loss Curve

The training loss curve is an effective way to understand the training dynamics of various methods. Following LISA, we present fine-tuning loss curves of LLaMa2-7B on the Alpaca-GPT4 dataset using Full FT, LoRA, LISA, and OWS in Figure 4-Right. At first glance, methods that directly fine-tune pre-trained weights (i.e., LISA and OWS) can better

Table 8: GSM8K scores/memory usage for fine-tuning LLaMA2-7B with various sampled layers γ . The results are averaged under three random seeds.

Method	Setting and Scores/Memory				
GaLore	r=8, $\gamma=32$	r=16, $\gamma=32$	r=32, $\gamma=32$	r=64, $\gamma=32$	r=128, $\gamma=32$
	19.1/35.6G	18.8/35.6G	18.4/35.8G	18.7/36.0G	18.2/36.5G
LISA	r=full, $\gamma=1$	r=full, $\gamma=2$	r=full, $\gamma=4$	r=full, $\gamma=8$	r=full, $\gamma=12$
	16.8/23G	18.8/25G	19.8/27G	19.9/32G	21.7/36G
OWS	r=128, $\gamma=1$	r=128, $\gamma=2$	r=128, $\gamma=4$	r=128, $\gamma=8$	r=128, $\gamma=12$
	20.0/21G	21.9/22G	23.5/23G	25.7/25G	27.8/27G

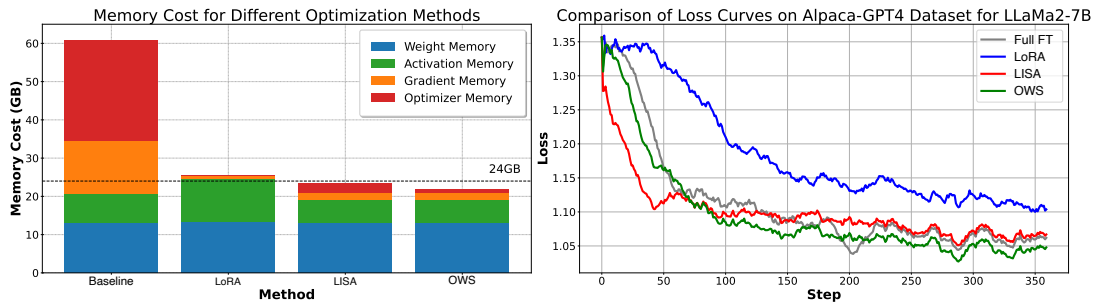


Figure 4: **Left:** Memory breakdown of various methods using LLaMa2-7B. **Right:** Fine-tuning loss of LLaMA2-7B on Alpaca GPT-4 dataset using various methods.

Table 9: Comparison with other layer-wise importance metrics, using LLaMA2-7B on Commonsense Reasoning.

Model	Sampling Method	Average
LlaMa2-7B	Uniform (Pan et al., 2024)	62.25
LlaMa2-7B	BI (Men et al., 2024)	62.15
LlaMa2-7B	RM (Samragh et al., 2023)	61.97
LlaMa2-7B	OWS (ours)	63.23

mimic the training landscape of full fine-tuning, compared to LoRA.

It is worth noting that while OWS initially falls short of LISA in the early phase of training, it gradually catches up after 60 iterations and eventually outperforms LISA with a lower loss. We conjecture that the underlying reason here is that the low-rank update of OWS is less accurate than the full-rank update of LISA at the beginning. However, as training progresses, OWS keeps updating the subspace, leading to an optimal one.

5 Related Work

Parameter-Efficient Fine-Tuning (PEFT). PEFT is proposed to reduce the prohibitive cost of LLM fine-tuning. Various techniques have been proposed in this dynamic field. For instance, prompt tuning only optimizes input tokens or embeddings while keeping the rest of the model frozen, as

demonstrated in studies (Lester et al., 2021; Li and Liang, 2021; Hambardzumyan et al., 2021; Zhong et al., 2021). Layer-freezing techniques (Liu et al., 2021b; Brock et al., 2017; Li et al., 2024) enhance training and fine-tuning efficiency by freezing parts of the layers. Adapter methods (Houlsby et al., 2019; He et al., 2021; Mahabadi et al., 2021; Diao et al., 2022), incorporate a small auxiliary module within the model’s architecture, which becomes the exclusive focus of updates during training, thus minimizing the number of trainable parameters and optimizer states. Among these techniques, Low-Rank Adaptation (LoRA) (Hu et al., 2021) gains massive attention by applying low-rank matrices to approximate weight changes during fine-tuning, which can be merged into the pre-trained weights, leading to no inference overhead. LoRA has been enhanced through various modifications (Zhang et al., 2023; Renduchintala et al., 2023; Sheng et al., 2023; Liu et al., 2024a; Kopiczko et al., 2023; Dettmers et al., 2024; Zhao et al., 2024) aimed at improving performance and efficiency. Recently, low-rank has also been explored to pre-train LLM from scratch (Lialin et al., 2023a; Zhao et al., 2024). GaLore (Zhao et al., 2024) projects the gradient into a low-rank subspace for the update to enable full-parameter learning while significantly reduc-

ing memory usage during optimization. BAdam (Luo et al., 2024) partitions the entire model into distinct blocks and utilizes a block coordinate descent framework to update each block individually, either in a deterministic or random sequence.

Layerwise Sampling for LLM Fine-tuning. Importance sampling is a powerful statistical technique used in machine learning to estimate properties of a particular distribution by sampling from a different, more convenient distribution. Recently, Pan et al. (2024) explored the idea of importance sampling to LLM fine-tuning, with the key idea of sampling only γ layers at each step to fine-tuning while keeping the rest of layers frozen. The proposed method, Layerwise Importance Sampled AdamW (LISA), outperforms LoRA by a large margin on various benchmarks and even outperforms full parameters training under certain settings. Inspired by LISA, our paper advances the performance of layerwise sampling for LLM fine-tuning, by addressing a couple of shortfalls of LISA.

6 Conclusion

In this paper, we study the sampling-based LLM fine-tuning, where at each iteration, only a few layers are sampled and fine-tuned, instead of the whole model. Specifically, we delve into recently-proposed LISA (Pan et al., 2024) and unveil two shortcomings that constrain its memory-performance trade-off: (1) The middle layers of LISA are sampled uniformly, which can result in suboptimal performance. (2) The sampled layers of LISA are fine-tuned in a full-rank manner, causing a significant memory increase as the number of sampled layers increases. To address these challenges, we introduced **OWS**, which assigns higher sampling probabilities to outlier-rich layers and incorporates low-rank gradient projection for improved memory efficiency. Our experiments on LLaMa2 and Mistral demonstrate that **OWS** significantly boosts performance while reducing memory usage compared to full-rank fine-tuning.

7 Limitations

While our proposed **OWS** approach demonstrates notable improvements in fine-tuning performance and memory efficiency, there are some factors to consider. Our evaluation primarily focuses on LLaMa2 and Mistral models, and further research could investigate the generalizability of **OWS** across a broader range of models and ar-

chitectures.

Acknowledgements

S. Liu is funded by the Royal Society with the Newton International Fellowship.

References

- Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. 2023. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*.
- Dan Biderman, Jose Gonzalez Ortiz, Jacob Portes, Mansheej Paul, Philip Greengard, Connor Jennings, Daniel King, Sam Havens, Vitaliy Chiley, Jonathan Frankle, et al. 2024. Lora learns less and forgets less. *arXiv preprint arXiv:2405.09673*.
- Andrew Brock, Theodore Lim, James M Ritchie, and Nick Weston. 2017. Freezeout: Accelerate training by progressively freezing layers. *arXiv preprint arXiv:1706.04983*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. 2022. Llm.int8(): 8-bit matrix multiplication for transformers at scale. *Advances in Neural Information Processing Systems (NeurIPS)*.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2024. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36.
- Shizhe Diao, Zhichao Huang, Ruijia Xu, Xuechun Li, Yong Lin, Xiao Zhou, and Tong Zhang. 2022. Black-box prompt learning for pre-trained language models. *arXiv preprint arXiv:2201.08531*.
- Karen Hambardzumyan, Hrant Khachatryan, and Jonathan May. 2021. Warp: Word-level adversarial reprogramming. *arXiv preprint arXiv:2101.00121*.
- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2021. Towards a unified view of parameter-efficient transfer learning. *arXiv preprint arXiv:2110.04366*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.

- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In *International conference on machine learning*, pages 2790–2799. PMLR.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Zhiqiang Hu, Lei Wang, Yihuai Lan, Wanyu Xu, Ee-Peng Lim, Lidong Bing, Xing Xu, Soujanya Poria, and Roy Ka-Wei Lee. 2023. Llm-adapters: An adapter family for parameter-efficient fine-tuning of large language models. *arXiv preprint arXiv:2304.01933*.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Xing Wang, Shuming Shi, and Zhaopeng Tu. 2023. Is chatgpt a good translator? yes with gpt-4 as the engine. *arXiv preprint arXiv:2301.08745*.
- Teun Kloek and Herman K Van Dijk. 1978. Bayesian estimates of equation system parameters: an application of integration by monte carlo. *Econometrica: Journal of the Econometric Society*, pages 1–19.
- Jan Kocoń, Igor Cichecki, Oliwier Kaszyca, Mateusz Kochanek, Dominika Szydło, Joanna Baran, Julita Bielaniec, Marcin Gruza, Arkadiusz Janz, Kamil Kanclerz, et al. 2023. Chatgpt: Jack of all trades, master of none. *Information Fusion*, 99:101861.
- Dawid Jan Kopiczko, Tijmen Blankevoort, and Yuki Markus Asano. 2023. Vera: Vector-based random matrix adaptation. *arXiv preprint arXiv:2310.11454*.
- Olga Kovaleva, Saurabh Kulshreshtha, Anna Rogers, and Anna Rumshisky. 2021. Bert busters: Outlier dimensions that disrupt transformers. *arXiv preprint arXiv:2105.06990*.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*.
- Sheng Li, Geng Yuan, Yue Dai, Youtao Zhang, Yanzhi Wang, and Xulong Tang. 2024. Smartfrz: An efficient training framework using attention-based layer freezing. *arXiv preprint arXiv:2401.16720*.
- Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*.
- Vladislav Lialin, Sherin Muckatira, Namrata Shiva-gunde, and Anna Rumshisky. 2023a. Relora: High-rank training through low-rank updates. In *Workshop on Advancing Neural Network Training: Computational Efficiency, Scalability, and Resource Optimization (WANT@ NeurIPS 2023)*.
- Vladislav Lialin, Namrata Shivagunde, Sherin Muckatira, and Anna Rumshisky. 2023b. Stack more layers differently: High-rank training through low-rank updates. *arXiv preprint arXiv:2307.05695*.
- Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. 2024a. Dora: Weight-decomposed low-rank adaptation. In *Forty-first International Conference on Machine Learning*.
- X Liu, Y Zheng, Z Du, M Ding, Y Qian, Z Yang, and J Tang. 2021a. Gpt understands, too. *arxiv. arXiv preprint arXiv:2103.10385*.
- Yongkang Liu, Yiqun Zhang, Qian Li, Tong Liu, Shi Feng, Daling Wang, Yifei Zhang, and Hinrich Schütze. 2024b. Hift: A hierarchical full parameter fine-tuning strategy. *arXiv preprint arXiv:2401.15207*.
- Yuhan Liu, Saurabh Agarwal, and Shivaram Venkataraman. 2021b. Autofreeze: Automatically freezing model blocks to accelerate fine-tuning. *arXiv preprint arXiv:2102.01386*.
- Haiquan Lu, Yefan Zhou, Shiwei Liu, Zhangyang Wang, Michael W Mahoney, and Yaoqing Yang. 2024. Al-phapruning: Using heavy-tailed self regularization theory for improved layer-wise pruning of large language models. *Advances in Neural Information Processing Systems*, 37:9117–9152.
- Qijun Luo, Hengxu Yu, and Xiao Li. 2024. Badam: A memory efficient full parameter training method for large language models. *arXiv preprint arXiv:2404.02827*.
- Rabeeh Karimi Mahabadi, Sebastian Ruder, Mostafa Dehghani, and James Henderson. 2021. Parameter-efficient multi-task fine-tuning for transformers via shared hypernetworks. *arXiv preprint arXiv:2106.04489*.
- Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D Lee, Danqi Chen, and Sanjeev Arora. 2023. Fine-tuning language models with just forward passes. *Advances in Neural Information Processing Systems*, 36:53038–53075.
- Xin Men, Mingyu Xu, Qingyu Zhang, Bingning Wang, Hongyu Lin, Yaojie Lu, Xianpei Han, and Weipeng Chen. 2024. Shortgpt: Layers in large language models are more redundant than you expect. *arXiv preprint arXiv:2403.03853*.
- Rui Pan, Xiang Liu, Shizhe Diao, Renjie Pi, Jipeng Zhang, Chi Han, and Tong Zhang. 2024. Lisa: Layer-wise importance sampling for memory-efficient

- large language model fine-tuning. *arXiv preprint arXiv:2403.17919*.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*.
- Giovanni Puccetti, Anna Rogers, Aleksandr Drozd, and Felice Dell’Orletta. 2022. Outliers dimensions that disrupt transformers are driven by frequency. *arXiv preprint arXiv:2205.11380*.
- Adithya Renduchintala, Tugrul Konuk, and Oleksii Kuchaiev. 2023. Tied-lora: Enhancing parameter efficiency of lora with weight tying. *arXiv preprint arXiv:2311.09578*.
- Mohammad Samragh, Mehrdad Farajtabar, Sachin Mehta, Raviteja Vemulapalli, Fartash Faghri, Devang Naik, Oncel Tuzel, and Mohammad Rastegari. 2023. Weight subcloning: direct initialization of transformers using larger pretrained ones. *arXiv preprint arXiv:2312.09299*.
- Ying Sheng, Shiyi Cao, Dacheng Li, Coleman Hooper, Nicholas Lee, Shuo Yang, Christopher Chou, Banghua Zhu, Lianmin Zheng, Kurt Keutzer, et al. 2023. S-lora: Serving thousands of concurrent lora adapters. *arXiv preprint arXiv:2311.03285*.
- Nigar M Shafiq Surameery and Mohammed Y Shakor. 2023. Use chat gpt to solve programming bugs. *International Journal of Information technology and Computer Engineering*, (31):17–22.
- Haoye Tian, Weiqi Lu, Tsz On Li, Xunzhu Tang, Shing-Chi Cheung, Jacques Klein, and Tegawendé F Bis-syandé. 2023. Is chatgpt the ultimate programming assistant—how far is it? *arXiv preprint arXiv:2304.11938*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Wenhan Xia, Chengwei Qin, and Elad Hazan. 2024. Chain of lora: Efficient fine-tuning of language models via residual learning. *arXiv preprint arXiv:2401.04151*.
- Lu Yin, You Wu, Zhenyu Zhang, Cheng-Yu Hsieh, Yaqing Wang, Yiling Jia, Mykola Pechenizkiy, Yi Liang, Zhangyang Wang, and Shiwei Liu. 2024. Outlier weighed layerwise sparsity (owl): A missing secret sauce for pruning llms to high sparsity. *In International Conference on Machine Learning*. PMLR.
- Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. 2023. Adaptive budget allocation for parameter-efficient fine-tuning. *In The Eleventh International Conference on Learning Representations*.
- Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong Tian. 2024. Galore: Memory-efficient llm training by gradient low-rank projection. *arXiv preprint arXiv:2403.03507*.
- Peilin Zhao and Tong Zhang. 2015. Stochastic optimization with importance sampling for regularized loss minimization. *In international conference on machine learning*, pages 1–9. PMLR.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36.
- Zexuan Zhong, Dan Friedman, and Danqi Chen. 2021. Factual probing is [mask]: Learning vs. learning to recall. *arXiv preprint arXiv:2104.05240*.

A OWS-Reverse

To further validate our approach, we introduce a new baseline: OWS-Reverse. This variant assigns lower sampling probabilities to layers with a higher proportion of outliers. As expected, OWS-Reverse performs the worst among the tested fine-tuning strategies, reinforcing our intuition about the importance of outlier-weighted prioritization in achieving better results.

Table 10: Comparison with varies baselines.

Method	MMLU	BoolQ	PIQA	SIQA	HellaSwag	WinoGrande	ARC-e	ARC-c	OBQA	Avg.
OWS-Reverse	49.4	81.9	77.8	33.4	59.1	80.1	79.3	50.2	38.2	61.0
Galore	49.6	81.8	79.4	32.9	60.7	79.6	79.8	49.4	37.6	61.2
LISA	49.6	82.0	79.9	33.5	59.7	79.6	80.4	51.1	38.8	61.6
OWS	52.6	85.4	80.7	34.2	60.3	82.2	80.6	51.0	39.1	62.9

B Hyperparameter Analysis

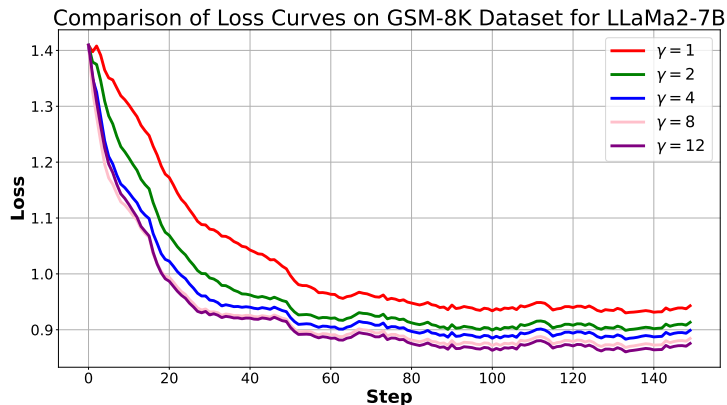


Figure 5: Fine-tuning loss of LLaMA2-7B using method OWS on the GSM-8K dataset with various sampled layers.

τ is the key hyperparameter to obtain the outlier ratio and sampling layers γ is also crucial to OWS To obtain intuitive and empirical guidance on these hyperparameter choices, we conduct ablation studies using LLaMA2-7B models with the GSM-8K dataset and report the results below.

Table 11: GSM scores for different τ values

Setting	$\tau = 3$	$\tau = 5$	$\tau = 7$	$\tau = 9$	$\tau = 11$	$\tau = 13$	$\tau = 15$	$\tau = 17$	$\tau = 19$
GSM Scores	19.18	19.41	20.04	20.62	21.15	20.24	20.17	20.47	19.79

We found that mid-range values of τ , such as 9, 11 and 13, generally lead to better performance. This may stem from the fact that the outliers screened by these values are more indicative of heavy-tailed properties. By default, we choose $\tau = 13$ for all experiments of OWS.

As for the sampling layer γ , it is not surprising that performance improves consistently with the sampling of more layers. OWS outperforms LISA with less memory usage across all sampling layer counts. This is attributed to OWS’s allocation of higher sampling probabilities to layers abundant in outliers, combined with its efficient low-rank gradient updating technique.

The training curve across different values of γ is depicted in Figure 5. Notably, fine-tuning with a higher γ leads to faster convergence and lower loss.

C Statistical Significance Test

We conducted experiments using 5 different seeds and reported the corresponding standard deviations. We do experiments with LISA and OWS to demonstrate the effectiveness of our proposed approach. For

MT-Bench, we provided the results evaluated using GPT-4o.

Table 12: Results of experiments for different models and methods evaluated on MT-Bench with 5 seeds and reported standard deviations.

Model	Method	BoolQ	PIQA	SIQA	HellaSwag	WinoGrande	ARC-e	ARC-c	OBQA
LLaMa2-7B	LISA	81.9 ± 0.22	79.6 ± 0.26	33.6 ± 0.11	59.6 ± 0.09	79.5 ± 0.16	80.3 ± 0.13	51.1 ± 0.11	39.2 ± 0.12
LLaMa2-7B	OWS	85.3 ± 0.19	80.8 ± 0.29	34.2 ± 0.14	60.2 ± 0.11	82.4 ± 0.18	80.8 ± 0.14	51.1 ± 0.12	39.6 ± 0.15
Mistral-7B	LISA	84.9 ± 0.21	82.7 ± 0.21	33.4 ± 0.11	64.4 ± 0.14	85.7 ± 0.16	83.6 ± 0.11	54.3 ± 0.10	40.5 ± 0.14
Mistral-7B	OWS	88.0 ± 0.24	84.0 ± 0.23	33.9 ± 0.11	66.4 ± 0.16	85.6 ± 0.09	84.1 ± 0.15	57.8 ± 0.14	40.5 ± 0.13

Additionally, we performed an independent samples t-test to assess the statistical significance of the performance difference between OWS and LISA. For example, in the LLaMa2-7B model, the t-test yields a t-statistic of -11.36 and a p-value of 3.41e-06, indicating that the performance improvements of OWS over LISA are statistically significant.

Table 13: Independent samples t-test results for the performance differences between OWS and LISA.

Model	t-statistic	p-value
LLaMa2-7B	-11.36	3.41e-06
Mistral-7B	-13.46	9.32e-07

D Training Configurations of OWS

We utilize Hugging Face and PyTorch for the implementation of our work.

Table 14: Hyperparameters used of OWS for fine-tuning LLaMa2-7B and Mistral-7B on the Commonsense Reasoning Benchmark.

Hyperparameter	LLaMa2-7B	Mistral-7B
Batch Size	16	16
Max. Sequence Length	512	512
Learning Rate	3e-4	3e-5
Scheduler	linear	linear
Training Epoch	1	1
Warmup Steps	0	0
dtype	bfloat16	bfloat16

Table 15: Hyperparameters used of OWS for fine-tuning LLaMa2-7B on various benchmarks.

Benchmarks	Commonsense Reasoning	MT-Bench	MMLU	GSM8K
Train Samples	170K	52K	99.8K	7.4K
Test Samples	22.4K	Alpaca-GPT4 (3.3K)	14K	1.3K
Batch Size	16	16	16	16
Max_length	512	512	512	512
Training Epoch	1	1	1	1
Learning Rate	3e-4	3e-4	3e-4	3e-4

E Pseudocode of GaLore

Following we present the pseudocode of OWS.

Algorithm 1: Outlier-Weighed Layerwise Sampling (OWS)

Require: number of layers N_L , number of training iterations T , sampling period K , sampled layers γ , rank level r , and $\mathcal{U}(0,1)$ refers to a uniform sampling.

% Before Training

for $\ell \leftarrow 1$ **to** N_L **do**

 Calculate outlier ratio D_j using the Equation 2

$$p_\ell \leftarrow \frac{\gamma D_\ell}{\sum_{j=1}^{N_L} D_j}$$

▷ *Mapping layerwise outlier distribution to sampling probability.*

% Training

for $i \leftarrow 0$ **to** $T/K - 1$ **do**

for $\ell \leftarrow 1$ **to** N_L **do**

if $\mathcal{U}(0,1) > p_\ell$ **then**

 Freeze layer ℓ

else

 Update the weights in layer ℓ

 grad = weight.grad

 lowrank_grad = **project**(grad)

 lowrank_update = **Adam_update**(lowrank_grad)

 update = **project_back**(lowrank_update)

 weight.data += update

▷ *OWS updates the in the low-rank subspace*

▷ *original space -> low-rank space*

▷ *update by Adam, Adafactor, etc.*

▷ *low-rank space -> original space*
