
4-bit Shampoo for Memory-Efficient Network Training

Sike Wang

Beijing Normal University
sikewang@mail.bnu.edu.cn

Pan Zhou

Singapore Management University
panzhou@smu.edu.sg

Jia Li[†]

Beijing Normal University
jiali@bnu.edu.cn

Hua Huang

Beijing Normal University
huahuang@bnu.edu.cn

Abstract

Second-order optimizers, maintaining a matrix termed a preconditioner, are superior to first-order optimizers in both theory and practice. The states forming the preconditioner and its inverse root restrict the maximum size of models trained by second-order optimizers. To address this, compressing 32-bit optimizer states to lower bitwidths has shown promise in reducing memory usage. However, current approaches only pertain to first-order optimizers. In this paper, we propose the first 4-bit second-order optimizers, exemplified by 4-bit Shampoo, maintaining performance similar to that of 32-bit ones. We show that quantizing the eigenvector matrix of the preconditioner in 4-bit Shampoo is remarkably better than quantizing the preconditioner itself both theoretically and experimentally. By rectifying the orthogonality of the quantized eigenvector matrix, we enhance the approximation of the preconditioner’s eigenvector matrix, which also benefits the computation of its inverse 4-th root. Besides, we find that linear square quantization slightly outperforms dynamic tree quantization when quantizing second-order optimizer states. Evaluation on various networks for image classification and natural language modeling demonstrates that our 4-bit Shampoo achieves comparable performance to its 32-bit counterpart while being more memory-efficient*.

1 Introduction

Deep neural networks (DNNs) have achieved great success in numerous fields, e.g., computer vision [20], natural language processing [38], and speech recognition [16]. A significant part of such success is attributed to first-order optimizers such as stochastic gradient descent with momentum (SGDM) [31] and AdamW [29]. Second-order optimizers, including K-FAC [30], Shampoo [18], AdaBK [41], CASPR [13], and Sophia [27], show great convergence properties, but often involve noticeable computation and memory costs. Anil et al. [2] provided several practical techniques for second-order optimizers to achieve substantial wall-clock time improvements over traditional first-order optimizers. The fast convergence property of second-order optimizers benefits from preconditioning the gradient with a matrix known as a preconditioner. The optimizer states for constructing the preconditioner and its inverse root can speed up optimization compared to first-order optimizers, but consume memory that could be used for model parameters, limiting the maximum model size trained within a given memory budget. With the increase in model size, the memory utilized by optimizer states can become a predominant factor in memory usage. This is the primary obstacle hindering the widespread use of second-order optimizers in the era of large models.

There are two main attempts to reduce memory consumed by optimizer states. Factorization uses low-rank approximation to optimizer states. This strategy has been applied to first-order optimizers [35, 3]

*Code is available at <https://github.com/Sike-Wang/low-bit-Shampoo>.

[†]Corresponding author.

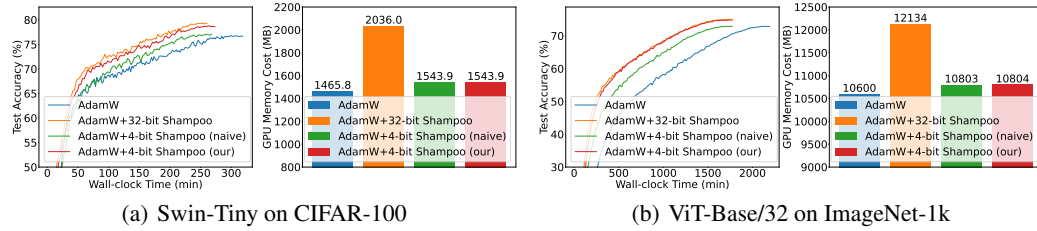


Figure 1: Visualization of test accuracies and total GPU memory costs of vision transformers. 4-bit Shampoo (naive) quantizes the preconditioner, while 4-bit Shampoo (our) quantizes its eigenvector matrix.

and second-order optimizers [14, 40]. In a comparable but distinct line of work, quantization utilizes low-bit to compress 32-bit optimizer states. Quantization is attractive due to its simplicity and wide applicability, which has been applied to first-order optimizers [8, 26]. Applying quantization to second-order optimizers poses a greater challenge, as first-order optimizers' states are elementwise, whereas second-order optimizers rely on matrix operations. To our knowledge, it has not been attempted before.

Contributions: In this paper, we present the first second-order optimizers with 4-bit optimizer states by taking Shampoo [18] as an example, while preserving the performance achieved with 32-bit optimizer states. While our focus is on Shampoo, we believe that our approach could also be applied to other second-order optimizers (see Table 4). Our main contributions are highlighted below.

Firstly, to maintain 32-bit performance, we propose quantizing the eigenvector matrix of a preconditioner in 4-bit Shampoo, rather than the preconditioner itself. The reason is that the small singular values of the preconditioner matter. Directly quantizing the preconditioner via block-wise quantization [8] at 4-bit precision can significantly alter the small singular values, leading to a drastic change in its inverse 4-th root and thus harming 4-bit Shampoo's performance. Quantizing the eigenvector matrix can help alleviate this issue, which is supported by experimental validation and theoretical insight. Additionally, with the eigenvector matrix, computing the inverse 4-th root is straightforward, ensuring that quantizing the eigenvector matrix does not lead to a rise in the total wall-clock time compared to quantizing the preconditioner (see Figure 1).

Secondly, we present two techniques for enhancing performance. As the eigenvector matrix of a preconditioner is orthogonal, we apply Björck orthonormalization [4] to rectify the orthogonality of the quantized eigenvector matrix, leading to improved approximation of preconditioner's eigenvector matrix and facilitating computation of its inverse 4-th root. Additionally, we observe that linear square quantization outperforms dynamic tree quantization [7] marginally when quantizing second-order optimizer states. The superiority of our developed 4-bit Shampoo is demonstrated in Figure 1.

Finally, we evaluate our 4-bit Shampoo on different image classification and natural language modeling tasks using convolutional neural network (CNN) and transformer architectures. Across all these benchmarks, our 4-bit Shampoo achieves similarly fast convergence comparable to its 32-bit counterpart, with no significant increase in losses for the trained models. Our 4-bit Shampoo uses less memory than its 32-bit counterpart, allowing for training of larger models with given resources.

2 Preliminaries

In this section, we present Shampoo and its implementation in our experiments. We also discuss quantization-based compression methods in a general formulation.

Notations. We use a non-bold letter like a or A to denote a scalar, a boldfaced lower-case letter like \mathbf{a} to denote a vector, and a boldfaced upper-case letter such as \mathbf{A} to denote a matrix. $\mathbf{u} = [u_i]^\top$ means that the i -th element of column vector \mathbf{u} is u_i and $\mathbf{U} = [\mathbf{u}_i]$ means the i -th column vector of matrix \mathbf{U} is \mathbf{u}_i . Let \mathbf{A} be a positive definite (PD) matrix and $s \in \mathbb{R}$, we define $\mathbf{A}^s = \mathbf{U} \mathbf{\Lambda}^s \mathbf{U}^\top$, where $\mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$ is the Singular Value Decomposition (SVD) of \mathbf{A} . $\text{tr}(\mathbf{A})$ represents the trace of a matrix \mathbf{A} . The inner product of two matrices \mathbf{A} and \mathbf{B} is denoted as $\langle \mathbf{A}, \mathbf{B} \rangle = \text{tr}(\mathbf{A}^\top \mathbf{B})$. The Frobenius norm of a matrix \mathbf{A} is $\|\mathbf{A}\|_F = \sqrt{\langle \mathbf{A}, \mathbf{A} \rangle}$. $\mathbf{A} \odot \mathbf{B}$ means the elementwise matrix product (Hadamard product).

$\text{Diag}(\mathbf{a})$ is a diagonal matrix with diagonal vector \mathbf{a} , while $\text{diag}(\mathbf{A})$ means the diagonal vector of matrix \mathbf{A} .

2.1 Shampoo for Matrices

The update rule of Shampoo in the matrix case combined with a first-order optimizer \mathcal{F} is

$$\text{Shampoo}(\mathbf{W}_{t-1}, \mathbf{L}_{t-1}, \mathbf{R}_{t-1}, \mathbf{s}_{t-1}, \mathbf{G}_t) = \begin{cases} \mathbf{L}_t = \mathbf{L}_{t-1} + \mathbf{G}_t \mathbf{G}_t^\top \\ \mathbf{R}_t = \mathbf{R}_{t-1} + \mathbf{G}_t^\top \mathbf{G}_t \\ \hat{\mathbf{G}}_t = \mathbf{L}_t^{-1/4} \mathbf{G}_t \mathbf{R}_t^{-1/4} \\ \tilde{\mathbf{G}}_t = \hat{\mathbf{G}}_t (\|\mathbf{G}_t\|_F / \|\hat{\mathbf{G}}_t\|_F) \\ \mathbf{W}_t, \mathbf{s}_t = \mathcal{F}(\mathbf{W}_{t-1}, \mathbf{s}_{t-1}, \tilde{\mathbf{G}}_t) \end{cases} \quad (1)$$

where \mathbf{W}_t is the model parameters in matrix form, \mathbf{L}_t and \mathbf{R}_t are called preconditioners, \mathbf{s}_t is the optimizer state of \mathcal{F} , and \mathbf{G}_t is the gradient at \mathbf{W}_{t-1} . Note that \mathbf{L}_t , \mathbf{R}_t , $\mathbf{L}_t^{-1/4}$, and $\mathbf{R}_t^{-1/4}$ are PD matrices. The penultimate step in (1) is the grafting trick [1], which enables Shampoo to roughly apply the well-tuned learning rate schedule of \mathcal{F} . The optimization variable \mathbf{W}_t does not represent all model parameters. It denotes a tensor of the model [18] or one block of a tensor [2]. In practice, we adopt an efficient and effective implementation of Shampoo for training DNNs following [2, 41] as described in Algorithm 4. In order to achieve efficient training, \mathbf{L}_t , \mathbf{R}_t , $\mathbf{L}_t^{-1/4}$, and $\mathbf{R}_t^{-1/4}$ are computed once every few hundred iterations. In this case, besides \mathbf{L}_t and \mathbf{R}_t , their inverse 4-th roots should also be stored in memory, as computing them is computationally expensive. So training large models with Shampoo can be memory-intensive, consuming a significant amount of memory.

2.2 Quantization-based Compression Methods

Quantizing updated optimizer states using a quantizer and then dequantizing them with a dequantizer prior to use is an effective method for conserving memory. We focus exclusively on vectors, as tensors can be reshaped into vectors.

Quantization. According to the idea in [8, 26], a b -bit quantizer \mathcal{Q} for p -dimensional real vectors is a mapping given by

$$\mathcal{Q} = (\mathcal{I} \circ \mathcal{N}, \mathcal{M}) : \mathbb{R}^p \rightarrow \mathbb{T}_b^p \times \mathbb{R}^p,$$

where \mathcal{N} is a normalization operator on \mathbb{R}^p , \mathcal{I} is an elementwise function mapping any real number to an element of $\mathbb{T}_b = \{0, 1, \dots, 2^b - 1\}$, and \mathcal{M} is a maximum operator on \mathbb{R}^p . For any $\mathbf{x} \in \mathbb{R}^p$, \mathcal{N} and \mathcal{M} satisfy $\mathcal{N}(\mathbf{x}) \odot \mathcal{M}(\mathbf{x}) = \mathbf{x}$.

A normalization operator \mathcal{N} for p -dimensional vectors is a transformation on \mathbb{R}^p . It scales each element of a vector $\mathbf{x} \in \mathbb{R}^p$ into $[-1, 1]$. A block-wise normalization operator for a p -dimensional vector $\mathbf{x} = [x_1, x_2, \dots, x_p]^\top$ is defined as

$$\mathcal{N}(\mathbf{x})_i = \frac{x_i}{\max_{j \in \mathbb{X}_i} \{x_j\}},$$

where $\mathcal{N}(\mathbf{x})_i$ is the i -th element of $\mathcal{N}(\mathbf{x})$, and \mathbb{X}_i is a set satisfying $i \in \mathbb{X}_i \subset \{1, \dots, p\}$. Usually, \mathbb{X}_i should also satisfy $\mathbb{X}_i = \mathbb{X}_j$ or $\mathbb{X}_i \cap \mathbb{X}_j = \emptyset$ for $i, j \in \{1, \dots, p\}$. In this case, for any $\mathbf{x} \in \mathbb{R}^p$, the number of different elements in $\mathcal{M}(\mathbf{x})$ is equal to the number of elements in set $\{\mathbb{X}_i | i = 1, \dots, p\}$. Meanwhile, the number of the elements in \mathbb{X}_i for any i should be as close as possible to a value called block size.

The mapping \mathcal{I} for $x \in \mathbb{R}$ in a b -bit quantizer \mathcal{Q} is defined as

$$\mathcal{I}(x) = \underset{j \in \mathbb{T}_b}{\text{argmin}} |x - \mathcal{R}(j)|,$$

where \mathcal{R} named quantization mapping is an elementwise function that maps any element in \mathbb{T}_b into $[-1, 1]$, and $|\cdot|$ is the absolute operator for a scalar. There are three typical quantization mappings: linear quantization, dynamic quantization, and quantile quantization. Their specifications and visualizations can be found in [8].

Dequantization. Given a b -bit quantizer $\mathcal{Q} = (\mathcal{I} \circ \mathcal{N}, \mathcal{M})$ for a p -dimensional real vector $\mathbf{x} \in \mathbb{R}^p$, the corresponding dequantizer \mathcal{D} is a mapping defined as

$$\mathcal{D}(\mathcal{Q}(\mathbf{x})) = \mathcal{D}(\mathcal{I} \circ \mathcal{N}(\mathbf{x}), \mathcal{M}(\mathbf{x})) = \mathcal{R}(\mathcal{I} \circ \mathcal{N}(\mathbf{x})) \odot \mathcal{M}(\mathbf{x}) : \mathbb{T}_b^p \times \mathbb{R}^p \rightarrow \mathbb{R}^p.$$

3 Methodology

In this section, we describe the design of our quantization-based compression method to realize 4-bit Shampoo with fast and high precision quantization. Let $\mathcal{Q}=(\mathcal{I} \circ \mathcal{N}, \mathcal{M})$ be a quantizer and \mathcal{D} be its corresponding dequantizer as described in Subsection 2.2.

3.1 Quantizing the Eigenvector Matrices

A naive approach to realize 4-bit Shampoo is applying the compression methods proposed in [8, 26] to \mathbf{L}_t , \mathbf{R}_t , $\mathbf{L}_t^{-1/4}$, and $\mathbf{R}_t^{-1/4}$ in Shampoo (see (1)). A slightly improved approach is to quantize the four PD matrices excluding their diagonal elements, which are typically much larger than their non-diagonal counterparts due to the non-negativity of the elements in $\text{diag}(\mathbf{G}_t \mathbf{G}_t^T)$ and $\text{diag}(\mathbf{G}_t^T \mathbf{G}_t)$.

However, the naive approach can cause large quantization errors at 4-bit precision. This is because the quantization errors (or called perturbations) of quantizing \mathbf{L}_t and \mathbf{R}_t will transfer to $\mathbf{L}_t^{-1/4}$ and $\mathbf{R}_t^{-1/4}$. To verify this, we first introduce two criteria to evaluate the quantization errors of matrices. We do not use the elementwise criterion in [8]. Let \mathbf{A} denote a 32-bit matrix, g represent a transformation (can formed by quantization), and f stand for a mapping, e.g., $f(\mathbf{A}) = \mathbf{A}^{-1/4}$. Then we define the normwise relative error (NRE) and angle error (AE) in f of g at \mathbf{A} as

$$\text{NRE} = \frac{\|f(\mathbf{A}) - f(g(\mathbf{A}))\|_F}{\|f(\mathbf{A})\|_F}, \quad \text{AE} = \arccos \left(\frac{\langle f(\mathbf{A}), f(g(\mathbf{A})) \rangle}{(\|f(\mathbf{A})\|_F \|f(g(\mathbf{A}))\|_F)} \right).$$

We choose two PD matrices of order 1200. The first one \mathbf{A}_1 is derived from the real world. It is a preconditioner in 32-bit Shampoo combined with AdamW for training a Swin-Tiny model. The second one $\mathbf{A}_2 = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$ is synthetic, constructed from a random orthogonal matrix \mathbf{U} and a diagonal matrix $\mathbf{\Lambda}$ with only two distinct diagonal values. Table 1 shows the quantization errors in $f(\mathbf{A}) = \mathbf{A}^{-1/4}$ of the naive approach at these two matrices, which are remarkably high. More analyses are given in Appendix D. The key point is that the singular values of $\mathbf{A}_i (i=1, 2)$ follow a specific distribution (see Figure 2). In this scenario, a slight perturbation of \mathbf{A}_i will significantly alter its small singular values, resulting in a drastic change to $\mathbf{A}_i^{-1/4}$.

To address this issue, we propose quantizing the eigenvector matrix of a preconditioner in Shampoo, rather than the preconditioner itself. Namely, a preconditioner \mathbf{A} is a PD matrix, and its SVD is $\mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$, where \mathbf{U} represents the eigenvector matrix and $\mathbf{\Lambda}$ denotes the singular value matrix. Given that $\mathbf{\Lambda}$ is a diagonal matrix, we can focus on quantizing \mathbf{U} using \mathcal{Q} while leaving $\mathbf{\Lambda}$ unchanged. From Table 1, one can observe that quantizing \mathbf{U} can significantly reduce the quantization errors. We will theoretically discuss the advantages of quantizing \mathbf{U} compared to quantizing \mathbf{A} in Section 4. In practice, the randomized SVD method [19] is adopted to compute the SVD of \mathbf{A} efficiently, as shown in [40]. We want to highlight that quantizing the original \mathbf{L}_t and \mathbf{R}_t in Shampoo involves significant computational burdens to compute their inverse 4-th roots $\mathbf{L}_t^{-1/4}$ and $\mathbf{R}_t^{-1/4}$, whereas quantizing the eigenvector matrices of \mathbf{L}_t and \mathbf{R}_t allows for rapid inverse root calculation. So the computational time required for both approaches is comparable (see Figure 1).

Table 1: Quantization errors in $\mathbf{A}^{-1/4}$ of different quantization schemes at a PD matrix \mathbf{A} . We employ block-wise normalization with a block size of 64. \mathbf{U} is the eigenvector matrix of \mathbf{A} , QM = quantized matrix, and OR = orthogonal rectification.

Real-world $\mathbf{A} = \mathbf{A}_1$						Synthetic $\mathbf{A} = \mathbf{A}_2$					
Mapping \mathcal{R}	Bit	QM	OR	NRE \downarrow	AE ($^\circ$) \downarrow	Mapping \mathcal{R}	Bit	QM	OR	NRE \downarrow	AE ($^\circ$) \downarrow
DT	8	\mathbf{A}	\times	0.2192	8.3014	DT	8	\mathbf{A}	\times	0.1896	10.877
	4	\mathbf{A}	\times	0.6241	17.319		4	\mathbf{A}	\times	0.4615	17.189
	4	\mathbf{U}	\times	0.0709	4.0426		4	\mathbf{U}	\times	0.1224	7.0144
	4	\mathbf{U}	\checkmark	0.0455	2.5615		4	\mathbf{U}	\checkmark	0.0878	4.9960
Linear-2	8	\mathbf{A}	\times	0.2164	7.9751	Linear-2	8	\mathbf{A}	\times	0.1310	7.4717
	4	\mathbf{A}	\times	0.6243	17.293		4	\mathbf{A}	\times	0.4465	15.338
	4	\mathbf{U}	\times	0.0543	3.1066		4	\mathbf{U}	\times	0.0942	5.3998
	4	\mathbf{U}	\checkmark	0.0343	1.9456		4	\mathbf{U}	\checkmark	0.0669	3.8166

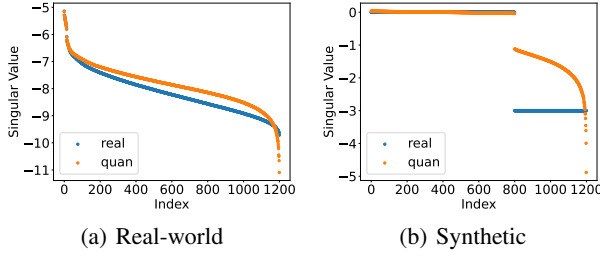


Figure 2: Singular value distributions of PD matrices (real) and their 4-bit compressions (quan) used in Table 1 with $\mathcal{R}=\text{DT}$, $\mathcal{Q}=\mathcal{M}$. Singular values are shown on a \log_{10} scale.

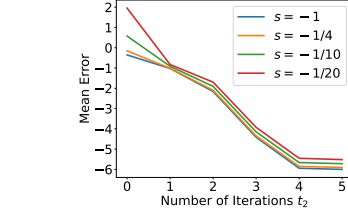


Figure 3: Elementwise mean errors between $(V_{t_2} \Lambda^s V_{t_2}^T)^{-1/s} (V_{t_2} \Lambda V_{t_2}^T)$ and identity matrix I . Mean errors are shown on a \log_{10} scale.

3.2 Rectifying the Orthogonality of Eigenvector Matrices

Let A be a PD matrix with SVD $U \Lambda U^T$. Note that the eigenvector matrix U is orthogonal, whereas $V = \mathcal{D}(\mathcal{Q}(U))$ may not be. To further mitigate the quantization errors mentioned in Subsection 3.1, we propose employing Björck orthonormalization [4] to orthogonalize V . Particularly, given $V_0 = V$, we iterate

$$V_t = 1.5V_{t-1} - 0.5V_{t-1}V_{t-1}^T V_{t-1}, \quad (2)$$

for $t_1 \geq 1$ times and take V_{t_1} as the rectified result. Equation (2) can also be interpreted as the gradient descent of problem $\min_V \|V^T V - I\|_F^2$ using a step size of 0.5, where I denotes the identity matrix. We empirically find that only one iteration (i.e., $t_1 = 1$) is enough. Table 1 illustrates the benefit of rectifying V into V_1 .

The update frequencies for the preconditioners and their inverse 4-th roots differ (see Algorithm 3). Given V and Λ , we also require orthogonal rectification to compute A^s rapidly for any $s \in \mathbb{R}$. The reason is as follows. It is easy to compute $A^s = U \Lambda^s U^T$ by definition. However, $U \Lambda^s U^T$ can be very sensitive to the orthogonality of U for $s < 0$, making $V \Lambda^s V^T$ largely deviate from $(V \Lambda V^T)^s \approx A^s$. Similarly, we can approximate A^s by $V_{t_2} \Lambda^s V_{t_2}^T$, where V_{t_2} is generated by (2). Figure 3 illustrates the elementwise mean errors between $(V_{t_2} \Lambda^s V_{t_2}^T)^{-1/s} (V_{t_2} \Lambda V_{t_2}^T)$ and I for various s and t_2 , where A is the real-world matrix used in Table 1. Based on the observation from Figure 3, we set $t_2 = 4$ in our experiments.

3.3 Selecting the Quantizer

The quantizer \mathcal{Q} is defined by the normalization operator \mathcal{N} and mapping \mathcal{R} , and \mathcal{N} is determined by \mathbb{X}_i . Since an eigenvector has a unit length, the elements in \mathbb{X}_i should belong to the same column of an eigenvector matrix, i.e., they are from the same eigenvector. Instead of employing dynamic tree (DT) quantization as mapping \mathcal{R} , we recommend utilizing linear square (Linear-2) quantization as \mathcal{R} , particularly when $b = 4$. Linear-2 quantization is defined as

$$\mathcal{R}(j) = \begin{cases} -(-1 + 2j/(2^b - 1))^2, & j < 2^{b-1} - 1; \\ 0, & j = 2^{b-1} - 1; \\ (-1 + 2j/(2^b - 1))^2, & j > 2^{b-1} - 1, \end{cases} \quad (3)$$

where $j \in \mathbb{T}_b = \{0, 1, \dots, 2^b - 1\}$. As shown in Table 1, Linear-2 quantization has lower quantization errors compared to DT quantization at 4-bit precision.

3.4 Overall Algorithm

We first describe the update processes of the preconditioners and their inverse 4-th roots in our 4-bit Shampoo. A preconditioner A is a PD matrix and its SVD is $U \Lambda U^T$. We can compress A into a pair $(\lambda, \bar{U}) = (\text{diag}(\Lambda), \mathcal{Q}(U))$ and decompress it into $(\Lambda, V) = (\text{Diag}(\lambda), \mathcal{D}(\bar{U}))$. Algorithm 1 (Preconditioner Update, PU) shows the update rule of A . Similarly, we compress $\hat{A} \approx A^{-1/4}$ into a pair $(a, \bar{A}) = (\text{diag}(\hat{A}), \mathcal{Q}(\hat{A} - \text{Diag}(a)))$ and decompress it into $\text{Diag}(a) + \mathcal{D}(\bar{A})$. Algorithm 2

(Preconditioner's Inverse 4-th Root Update, PIRU) gives the update rule of $\hat{\mathbf{A}}$. Based on the above update rules, we can summarize our 4-bit Shampoo in Algorithm 3. Note that we omit some input parameters of PU and PIRU because they can be found in Algorithm 3 in the same form.

Algorithm 1 PU($\lambda, \bar{\mathbf{U}}, \mathbf{M}$) Input: singular value vector λ , quantized eigenvector matrix $\bar{\mathbf{U}}$, \mathbf{M} , number of iterations t_1 for rectification, exponential decay rate $\beta \in (0, 1)$, \mathcal{Q} and \mathcal{D} 1: $\mathbf{\Lambda} = \text{Diag}(\lambda)$, $\mathbf{V} = \mathcal{D}(\bar{\mathbf{U}})$ 2: Rectify \mathbf{V} by iterating (2) t_1 times 3: $\mathbf{A} = \beta \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top + (1 - \beta) \mathbf{M}$ 4: Compute $\mathbf{A} = \mathbf{P} \mathbf{\Sigma} \mathbf{P}^\top$ by randomized SVD 5: return $\text{diag}(\mathbf{\Sigma})$, $\mathcal{Q}(\mathbf{P})$	Algorithm 2 PIRU($\lambda, \bar{\mathbf{U}}$) Input: singular value vector λ , quantized eigenvector matrix $\bar{\mathbf{U}}$, number of iterations t_2 for rectification, dampening term $\epsilon \mathbf{I}$, \mathcal{Q} and \mathcal{D} 1: $\mathbf{\Lambda} = \text{Diag}(\lambda)$, $\mathbf{V} = \mathcal{D}(\bar{\mathbf{U}})$ 2: Rectify \mathbf{V} by iterating (2) t_2 times 3: $\hat{\mathbf{A}} = \mathbf{V}(\mathbf{\Lambda} + \max\{\lambda\} \epsilon \mathbf{I})^{-1/4} \mathbf{V}^\top$ 4: $\mathbf{a} = \text{diag}(\hat{\mathbf{A}})$ 5: return \mathbf{a} , $\mathcal{Q}(\hat{\mathbf{A}} - \text{Diag}(\mathbf{a}))$
Algorithm 3 Practical 4-bit Shampoo Input: $\mathbf{W}_0 \in \mathbb{R}^{m \times n}$, $\mathbf{L}_0 = \epsilon \mathbf{I}_m$, $\mathbf{R}_0 = \epsilon \mathbf{I}_n$, $\hat{\mathbf{L}}_0 = \mathbf{I}_m$, $\hat{\mathbf{R}}_0 = \mathbf{I}_n$, $\beta \in (0, 1)$, t_1, t_2 , update interval T_1 , update interval T_2 , total number of steps T , first-order optimizer \mathcal{F} , first-order optimizer state $\mathbf{s}_0 = \mathbf{0}$, 4-bit quantizer \mathcal{Q} and its corresponding dequantizer \mathcal{D} . Output: final parameter \mathbf{W}_T . 1: $\lambda_{0,L} = \text{diag}(\mathbf{L}_0)$, $\bar{\mathbf{U}}_{0,L} = \mathcal{Q}(\mathbf{I}_m)$; $\lambda_{0,R} = \text{diag}(\mathbf{R}_0)$, $\bar{\mathbf{U}}_{0,R} = \mathcal{Q}(\mathbf{I}_n)$ 2: $\mathbf{l}_0 = \text{diag}(\hat{\mathbf{L}}_0)$, $\bar{\mathbf{L}}_0 = \mathcal{Q}(\mathbf{0})$; $\mathbf{r}_0 = \text{diag}(\hat{\mathbf{R}}_0)$, $\bar{\mathbf{R}}_0 = \mathcal{Q}(\mathbf{0})$ 3: for $t = 1, 2, \dots, T$ do 4: Receive loss function $\mathcal{L}_t : \mathbb{R}^{m \times n} \mapsto \mathbb{R}$ and compute gradient $\mathbf{G}_t = \nabla \mathcal{L}_t(\mathbf{W}_t)$ 5: if $t \% T_1 \equiv 0$ then 6: $\lambda_{t,L}, \bar{\mathbf{U}}_{t,L} = \text{PU}(\lambda_{t-1,L}, \bar{\mathbf{U}}_{t-1,L}, \mathbf{G}_t \mathbf{G}_t^\top)$; $\lambda_{t,R}, \bar{\mathbf{U}}_{t,R} = \text{PU}(\lambda_{t-1,R}, \bar{\mathbf{U}}_{t-1,R}, \mathbf{G}_t^\top \mathbf{G}_t)$ 7: else 8: $\lambda_{t,L}, \bar{\mathbf{U}}_{t,L} = \lambda_{t-1,L}, \bar{\mathbf{U}}_{t-1,L}$; $\lambda_{t,R}, \bar{\mathbf{U}}_{t,R} = \lambda_{t-1,R}, \bar{\mathbf{U}}_{t-1,R}$ 9: if $t \% T_2 \equiv 0$ then 10: $\mathbf{l}_t, \bar{\mathbf{L}}_t = \text{PIRU}(\lambda_{t,L}, \bar{\mathbf{U}}_{t,L})$; $\mathbf{r}_t, \bar{\mathbf{R}}_t = \text{PIRU}(\lambda_{t,R}, \bar{\mathbf{U}}_{t,R})$ 11: else 12: $\mathbf{l}_t, \bar{\mathbf{L}}_t = \mathbf{l}_{t-1}, \bar{\mathbf{L}}_{t-1}$; $\mathbf{r}_t, \bar{\mathbf{R}}_t = \mathbf{r}_{t-1}, \bar{\mathbf{R}}_{t-1}$ 13: $\hat{\mathbf{L}}_t = \text{Diag}(\mathbf{l}_t) + \mathcal{D}(\bar{\mathbf{L}}_t)$; $\hat{\mathbf{R}}_t = \text{Diag}(\mathbf{r}_t) + \mathcal{D}(\bar{\mathbf{R}}_t)$ 14: $\hat{\mathbf{G}}_t = \hat{\mathbf{L}}_t \mathbf{G}_t \hat{\mathbf{R}}_t$; $\tilde{\mathbf{G}}_t = \hat{\mathbf{G}}_t (\ \mathbf{G}_t\ _F / \ \hat{\mathbf{G}}_t\ _F)$ 15: $\mathbf{W}_t, \mathbf{s}_t = \mathcal{F}(\mathbf{W}_{t-1}, \mathbf{s}_{t-1}, \tilde{\mathbf{G}}_t)$	

4 Theoretical Analysis

In this section, we analyze why quantizing the eigenvector matrix of a preconditioner in Shampoo is better than quantizing the preconditioner itself under a certain singular value distribution. Furthermore, we consider quantization as a perturbation and prove the convergence of the perturbed Shampoo (Algorithm 6) in Appendix E. The following lemma reveals some good properties of perturbing the eigenvector matrix of a PD matrix.

Lemma 1. Let \mathbf{A} be a PD matrix whose SVD is $\mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$, where $\mathbf{U} = [\mathbf{u}_i]$ is an orthogonal matrix and $\mathbf{\Lambda} = \text{diag}([\lambda_i]^\top)$ is a diagonal matrix. Given a perturbation $\Delta \mathbf{U} = [\Delta \mathbf{u}_i]$ and $s \in \mathbb{R}$, we define $\mathbf{B} := (\mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top)^s$ and $\Delta \mathbf{B} := ((\mathbf{U} + \Delta \mathbf{U}) \mathbf{\Lambda} (\mathbf{U} + \Delta \mathbf{U})^\top)^s - \mathbf{B}$.

(1) If $\mathbf{U} + \Delta \mathbf{U}$ is orthogonal and there exists $\alpha \in \mathbb{R}$ such that $\|\Delta \mathbf{u}_i\|_2 \leq \alpha$, then

$$\frac{\|\Delta \mathbf{B}\|_F}{\|\mathbf{B}\|_F} \leq 2\alpha.$$

(2) If $\mathbf{U} + \Delta \mathbf{U}$ is orthogonal and there exists $\beta \in \mathbb{R}$ such that $\langle \mathbf{u}_i, \mathbf{u}_i + \Delta \mathbf{u}_i \rangle \geq 1 - \beta \geq 0$, then

$$\frac{\langle \mathbf{B}, \mathbf{B} + \Delta \mathbf{B} \rangle}{\|\mathbf{B}\|_F \|\mathbf{B} + \Delta \mathbf{B}\|_F} \geq (1 - \beta)^2.$$

From Lemma 1, it is evident that the normwise relative error and angle error in $f(\mathbf{A}) = \mathbf{A}^s$ of perturbing \mathbf{U} at $\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$ are independent of $\mathbf{\Lambda}$ and s . Moreover, these errors are well-bounded under some mild conditions. Empirically, for 4-bit quantization, $\alpha = 0.1$ and $\beta = 0.005$ roughly meet the conditions of Lemma 1, leading to $\frac{\|\Delta\mathbf{B}\|_F}{\|\mathbf{B}\|_F} \leq 0.2$ and $\frac{\langle \mathbf{B}, \mathbf{B} + \Delta\mathbf{B} \rangle}{\|\mathbf{B}\|_F \|\mathbf{B} + \Delta\mathbf{B}\|_F} \geq 0.99$.

It is very complicated to generally analyze the perturbation in $f(\mathbf{A}) = \mathbf{A}^s$ of perturbing \mathbf{A} . Thus, we focus on perturbing the singular values of \mathbf{A} . For simplicity, we assume that both \mathbf{A} and $\mathbf{A} + \Delta\mathbf{A}$ have only two distinct singular values, where $\Delta\mathbf{A}$ is a perturbation of \mathbf{A} . The following lemma gives the perturbation in \mathbf{A}^s of perturbing the smaller singular value of \mathbf{A} .

Lemma 2. Let \mathbf{A} be a PD matrix of order $m+n$ whose SVD is $\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$, where $m, n \in \mathbb{N}_+$, $n = lm$, $\mathbf{U} = [\mathbf{u}_i]$ is an orthogonal matrix and $\mathbf{\Lambda} = \text{diag}([\lambda_i]^\top)$ is a diagonal matrix. Assume that $\mathbf{\Lambda} = \text{diag}([c\lambda\mathbf{1}_{m \times 1}^\top, \lambda\mathbf{1}_{n \times 1}^\top]^\top)$, $c \geq 1$, and $\lambda > 0$. Given a perturbation $\Delta\mathbf{\Lambda} = \text{diag}([\mathbf{0}_{m \times 1}^\top, \Delta\lambda_{n \times 1}^\top]^\top)$ and $s \in \mathbb{R}$, we define $\mathbf{B} := (\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top)^s$ and $\Delta\mathbf{B} := (\mathbf{U}(\mathbf{\Lambda} + \Delta\mathbf{\Lambda})\mathbf{U}^\top)^s - \mathbf{B}$.

(1) If $\Delta\lambda_{n \times 1} = (k-1)\lambda\mathbf{1}_{n \times 1}$ where $k > 0$, then

$$\frac{\|\Delta\mathbf{B}\|_F}{\|\mathbf{B}\|_F} = \frac{\sqrt{l}|k^s - 1|}{\sqrt{c^{2s} + l}} = h_1(s, l).$$

Moreover, $h_1(s, l)$ decreases monotonically with s over $(-\infty, 0)$ and increases monotonically with l over $(0, +\infty)$.

(2) If $\Delta\lambda_{n \times 1} = (tc-1)\lambda\mathbf{1}_{n \times 1}$ where $t > 0$, then

$$\frac{\langle \mathbf{B}, \mathbf{B} + \Delta\mathbf{B} \rangle}{\|\mathbf{B}\|_F \|\mathbf{B} + \Delta\mathbf{B}\|_F} = \frac{lt^s + c^s}{\sqrt{(1 + lt^{2s})(l + c^{2s})}} = h_2(l).$$

Moreover, $h_2(l)$ decreases monotonically with l over $(0, (c/t)^s]$ and increases monotonically with l over $((c/t)^s, +\infty)$.

(3) If $\Delta\lambda_{n \times 1} = (tc-1)\lambda\mathbf{1}_{n \times 1}$ where $k = tc > 0$ and $l = (c/t)^s$, then

$$\frac{\|\Delta\mathbf{B}\|_F}{\|\mathbf{B}\|_F} = \frac{|k^s - 1|}{\sqrt{k^s + 1}}, \quad \frac{\langle \mathbf{B}, \mathbf{B} + \Delta\mathbf{B} \rangle}{\|\mathbf{B}\|_F \|\mathbf{B} + \Delta\mathbf{B}\|_F} = \frac{2}{\sqrt{2 + k^s + 1/k^s}}.$$

Let us make some comments on the above lemma. First, from Lemma 2(1) we have $h_1(1, l) = \frac{\|\Delta\mathbf{A}\|_F}{\|\mathbf{A}\|_F} = \frac{\sqrt{l}|k-1|}{\sqrt{c^{2s} + l}}$. If $k \geq 1$, $\frac{\|\Delta\mathbf{A}\|_F}{\|\mathbf{A}\|_F} = \frac{\|\Delta\mathbf{\Lambda}\|_F}{\|\mathbf{\Lambda}\|_F}$ is bounded by $\frac{k}{c}\sqrt{l} = t\sqrt{l}$. Second, if $k = tc \geq 1$ and $s < 0$, one can deduce $h_2(l) \geq \sqrt{lt^{2s}/(1 + lt^{2s})}$ from Lemma 2(2), which indicates that a small lt^{2s} is needed to achieve small $h_2(l)$. We can set $t = 0.02$ to simulate 4-bit quantization. Based on Lemma 1 and Lemma 2(3), we have the following proposition.

Proposition 1. Let \mathbf{A} be a PD matrix of order $m+n$ whose SVD is $\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$, where $m, n \in \mathbb{N}_+$, $n = lm$, $\mathbf{U} = [\mathbf{u}_i]$ is an orthogonal matrix, $\mathbf{\Lambda} = \text{diag}([c\lambda\mathbf{1}_{m \times 1}^\top, \lambda\mathbf{1}_{n \times 1}^\top]^\top)$, $c \geq 1000$, and $\lambda > 0$. Given $\Delta\mathbf{U} = [\Delta\mathbf{u}_i]$, $\Delta\mathbf{\Lambda} = \text{diag}([\mathbf{0}_{m \times 1}^\top, \Delta\lambda_{n \times 1}^\top]^\top)$, and $s \leq -0.25$, we define $\mathbf{B} := (\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top)^s$, $\mathbf{B}_1 := ((\mathbf{U} + \Delta\mathbf{U})\mathbf{\Lambda}(\mathbf{U} + \Delta\mathbf{U})^\top)^s$, and $\mathbf{B}_2 := (\mathbf{U}(\mathbf{\Lambda} + \Delta\mathbf{\Lambda})\mathbf{U}^\top)^s$. If $\mathbf{U} + \Delta\mathbf{U}$ is orthogonal, $\|\Delta\mathbf{u}_i\|_2 \leq 0.1$, $\langle \mathbf{u}_i, \Delta\mathbf{u}_i \rangle \geq -0.005$, $\Delta\lambda_{n \times 1} = (0.02c-1)\lambda\mathbf{1}_{n \times 1}$, and $l = (c/0.02)^s$, then

$$2 \frac{\|\mathbf{B}_1 - \mathbf{B}\|_F}{\|\mathbf{B}\|_F} \leq 0.4 \leq \frac{\|\mathbf{B}_2 - \mathbf{B}\|_F}{\|\mathbf{B}\|_F}, \quad 6 \left(1 - \frac{\langle \mathbf{B}, \mathbf{B}_1 \rangle}{\|\mathbf{B}\|_F \|\mathbf{B}_1\|_F} \right) \leq 0.06 \leq \left(1 - \frac{\langle \mathbf{B}, \mathbf{B}_2 \rangle}{\|\mathbf{B}\|_F \|\mathbf{B}_2\|_F} \right).$$

Proposition 1 requires very strong assumptions. Nevertheless, it provides insight into why quantizing \mathbf{A} can result in a greater normwise relative error and angle error in \mathbf{A}^s , compared to quantizing \mathbf{U} . Complete proofs of Lemma 1, Lemma 2, and Proposition 1 can be found in Appendix F.

5 Experiments

In this section, we compare our 4-bit Shampoo combined with SGDM or AdamW to their 32-bit counterparts, as well as the first-order optimizers on various image classification tasks. See more experimental results on image classification and natural language modeling tasks in Appendix H.

Models, datasets, and hyperparameters. We train VGG19 [36], ResNet34 [20], ViT-Small [10], and Swin-Tiny [28] on the CIFAR-100 [23] and Tiny-ImageNet [24] datasets with one RTX3060Ti GPU, and train ResNet50 and ViT-Base/32 on the ImageNet-1k dataset [34] with one A800 GPU.

Table 2: Performance, wall-clock time and memory cost on various image classification tasks. TA = test accuracy, WCT = wall-clock time, and TMC = total GPU memory cost.

Dataset	Model	Optimizer	TA (%)	WCT (min)	TMC (MB)
CIFAR-100	VGG19	SGDM	74.14	97.70	512.17
		SGDM + 32-bit Shampoo	74.54	84.45	979.13
		SGDM + 4-bit Shampoo	74.74	92.51	577.14
	ResNet34	SGDM	78.98	170.1	822.03
		SGDM + 32-bit Shampoo	79.71	147.2	1441.8
		SGDM + 4-bit Shampoo	79.17	155.8	908.40
	ViT-Small	AdamW	74.34	668.1	2720.0
		AdamW + 32-bit Shampoo	77.50	498.7	3252.0
		AdamW + 4-bit Shampoo	77.22	510.8	2791.7
	Swin-Tiny	AdamW	76.69	318.6	1465.8
		AdamW + 32-bit Shampoo	79.34	260.8	2036.0
		AdamW + 4-bit Shampoo	78.63	273.3	1543.9
Tiny-ImageNet	VGG19	SGDM	61.53	172.0	1062.3
		SGDM + 32-bit Shampoo	63.39	136.5	1531.9
		SGDM + 4-bit Shampoo	62.84	143.8	1127.3
	ResNet34	SGDM	67.10	432.1	2304.0
		SGDM + 32-bit Shampoo	67.90	313.0	2924.3
		SGDM + 4-bit Shampoo	67.95	329.3	2390.4
	ViT-Small	AdamW	54.66	1274	2730.1
		AdamW + 32-bit Shampoo	57.11	953.9	3261.1
		AdamW + 4-bit Shampoo	57.15	970.3	2801.9
	Swin-Tiny	AdamW	58.77	701.9	1789.9
		AdamW + 32-bit Shampoo	61.74	565.3	2362.8
		AdamW + 4-bit Shampoo	62.24	582.7	1868.1
ImageNet-1k	ResNet50	SGDM	76.70	2134	11307
		SGDM + 32-bit Shampoo	77.07	1910	11937
		SGDM + 4-bit Shampoo	76.92	1970	11396
	ViT-Base/32	AdamW	72.87	2190	10600
		AdamW + 32-bit Shampoo	75.03	1774	12134
		AdamW + 4-bit Shampoo	74.78	1770	10804

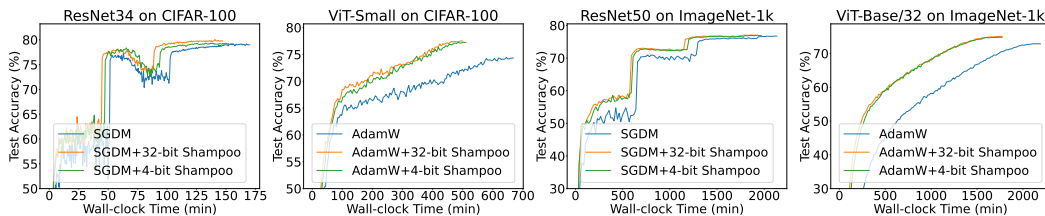


Figure 4: Visualization of test accuracies on the CIFAR-100 and ImageNet-1k datasets.

For hyperparameter settings, we mainly follow [41] to train CNNs and [25, 44] to train vision transformers. For all the tasks, we keep the common hyperparameters of optimizers the same values. See Appendix G for experimental details.

Main results. We show the performance, wall-clock time, and memory cost in Table 2. First-order optimizers run 1.2x to 1.5x epochs, resulting in longer wall-clock time, yet yielding lower test accuracies compared to second-order optimizers. In comparison to 32-bit Shampoo, our 4-bit Shampoo shows comparable test accuracies with differences ranging from -0.7% to 0.5%, increases in wall-clock time varying from -0.2% to 9.5%, and memory savings of 4.5% to 41%. Compared to the first-order optimizers, the memory costs of our 4-bit Shampoo only rise by 0.8% to 12.7%. This represents a significant advancement in the utilization of second-order optimizers. Following [26], we report the total peak GPU memory consumption rather than the optimizer’s peak GPU memory consumption. Our main focus is on quantizing the states for constructing preconditioners and their inverse roots, which are approximately 7x smaller for 4-bit Shampoo compared to 32-bit Shampoo (see Appendix G). Figure 4 shows the test accuracy curves on the CIFAR-100 and ImageNet-1k

Table 3: Ablation study on the impact of different quantization techniques to Swin-Tiny training on the CIFAR-100 dataset. U is the eigenvector matrix of a preconditioner A . QM = quantized matrix, OR = orthogonal rectification in Algorithm 1, TL = training loss, and TA = test accuracy.

4-bit					3-bit				
Mapping \mathcal{R}	QM	OR	TL	TA (%)	Mapping \mathcal{R}	QM	OR	TL	TA (%)
Linear-2	A	\times	1.631	76.95	Linear-2	A	\times	1.648	76.70
DT	U	\times	1.569	78.70	DT	U	\times	NaN	-
Linear-2	U	\times	1.566	78.22	Linear-2	U	\times	NaN	-
Linear-2	U	\checkmark	1.551	78.63	Linear-2	U	\checkmark	1.572	78.53

datasets. The test accuracy curves of 4-bit Shampoo and 32-bit Shampoo are very close, both of which are above the test accuracy curves of the first-order optimizers.

Ablations. We investigate the effectiveness of our proposed quantization techniques. Table 3 indicates that quantizing the eigenvector matrix of a preconditioner is crucial for b -bit ($b = 3, 4$) Shampoo to maintain 32-bit performance, and orthogonal rectification is highly beneficial for 3-bit Shampoo. As for quantization mapping, linear square (Linear-2) quantization is comparable to dynamic tree (DT) quantization. We further apply our 4-bit quantization techniques to K-FAC [30], AdaBK [41] and CASPR [13] and the results are shown in Table 4. We can see that the 4-bit optimizers match the performance of their 32-bit counterparts, and reduce memory by over 20%.

6 Related Work

Second-order optimizers. Different second-order optimizers apply different second-order information. Hessian-based optimizers [39, 27] use the Hessian matrix or its approximation. Fisher-based optimizers [30, 41] utilize the covariance matrix of the accumulated gradients or its approximation based on Kronecker product. Shampoo [18] and CASPR [13] approximate the full AdaGrad [12] preconditioner by a set of small preconditioning matrices.

Memory efficient optimizers based on factorization. Adafactor [35] employs the outer product of two vectors to approximate the second moment of Adam [22]. SM3 [3] considers approximating the second moment of Adam by its covers' statistics. [14] and [40] reduce memory cost of the preconditioner in a second-order optimizer with its low-rank approximation through truncated SVD.

Memory efficient optimizers based on quantization. Dettmers et al. [8] introduce block-wise dynamic quantization that enables the use of first-order optimizers with 8-bit states. Li et al. [26] push the optimizer states of Adam/AdamW to 4-bit.

Table 4: Performance and memory cost of training Swin-Tiny on CIFAR-100. TA = test accuracy and TMC = total GPU memory cost.

Optimizer	TA (%)	TMC (MB)
AdamW+32-bit K-FAC	78.20	2388.0
AdamW+4-bit K-FAC	78.56	1878.3
32-bit AdamW_BK	79.28	2388.0
4-bit AdamW_BK	79.34	1878.3
AdamW+32-bit CASPR	78.82	2034.6
AdamW+4-bit CASPR	78.80	1543.9

7 Conclusions, Limitations, and Broader Impact

We propose 4-bit Shampoo, the first low-bit second-order optimizer, designed for memory-efficient training of DNNs. We find that quantizing the eigenvector matrix of the preconditioner is essential to minimize quantization errors in its inverse 4-th root at 4-bit precision, given its sensitivity to alterations in small singular values. We further introduce orthogonal rectification and linear square quantization mapping to improve performance. 4-bit Shampoo achieves lossless performance to 32-bit counterpart in training different DNNs on various tasks.

Limitations. Preconditioners in Shampoo are symmetric matrices and can be stored as upper triangular matrices, saving almost half of the memory usage. However, the eigenvector matrix of a preconditioner is not symmetric, causing an 8-bit preconditioner to occupy the same memory as its 4-bit eigenvector matrix. Notably, a comparison of Table 1 and Table 7 in Appendix D shows that the 4-bit quantization of the eigenvector matrix has smaller quantization errors than the 8-bit quantization of the preconditioner. Our evaluation is currently limited to image classification and natural language modeling tasks. Due to limitations in computing resources, we do not test our 4-bit Shampoo on large-scale models with billions of parameters.

Broader Impact. Our work can facilitate training large models with second-order optimizers. This could open up new research possibilities that were previously unattainable due to GPU memory constraints, especially benefiting researchers with limited resources.

Acknowledgments and Disclosure of Funding

Jia Li and Hua Huang were supported by the NSF of China (grant no. 62131003). Jia Li was also supported by the NSF of China (grant no. 62102034). Pan Zhou was supported by the Singapore Ministry of Education (MOE) Academic Research Fund (AcRF) Tier 1 grants (project ID: 23-SIS-SMU-028 and 23-SIS-SMU-070).

References

- [1] Naman Agarwal, Rohan Anil, Elad Hazan, Tomer Koren, and Cyril Zhang. Disentangling adaptive gradient methods from learning rates. *arXiv preprint arXiv:2002.11803*, 2020.
- [2] Rohan Anil, Vineet Gupta, Tomer Koren, Kevin Regan, and Yoram Singer. Scalable second order optimization for deep learning. *arXiv preprint arXiv:2002.09018*, 2020.
- [3] Rohan Anil, Vineet Gupta, Tomer Koren, and Yoram Singer. Memory efficient adaptive optimization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [4] Å. Björck and C. Bowie. An iterative algorithm for computing the best estimate of an orthogonal matrix. *SIAM Journal on Numerical Analysis*, 8(2):358–364, 1971.
- [5] R.L. Burden, J.D. Faires, and A.M. Burden. *Numerical Analysis*. Cengage Learning, 2015.
- [6] Aaron Defazio, Xingyu Yang, Harsh Mehta, Konstantin Mishchenko, Ahmed Khaled, and Ashok Cutkosky. The road less scheduled. *arXiv preprint arXiv:2405.15682*, 2024.
- [7] Tim Dettmers. 8-bit approximations for parallelism in deep learning. In *Proceedings of the International Conference on Learning Representations*, 2016.
- [8] Tim Dettmers, Mike Lewis, Sam Shleifer, and Luke Zettlemoyer. 8-bit optimizers via block-wise quantization. In *Proceedings of the International Conference on Learning Representations*, 2022.
- [9] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. QLoRA: Efficient finetuning of quantized LLMs. *Advances in Neural Information Processing Systems*, 2023.
- [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *Proceedings of the International Conference on Learning Representations*, 2021.
- [11] Timothy Dozat. Incorporating Nesterov momentum into Adam. In *Proceedings of the International Conference on Learning Representations Workshop*, 2016.
- [12] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(61):2121–2159, 2011.
- [13] Sai Surya Duvvuri, Fnu Devvrit, Rohan Anil, Cho-Jui Hsieh, and Inderjit S Dhillon. Combining axes preconditioners through Kronecker approximation for deep learning. In *Proceedings of the International Conference on Learning Representations*, 2024.
- [14] Vladimir Feinberg, Xinyi Chen, Y. Jennifer Sun, Rohan Anil, and Elad Hazan. Sketchy: Memory-efficient adaptive regularization with frequent directions. *Advances in Neural Information Processing Systems*, 2023.
- [15] Elias Frantar, Eldar Kurtic, and Dan Alistarh. M-FAC: Efficient matrix-free approximations of second-order information. *Advances in Neural Information Processing Systems*, 2021.

- [16] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, and Ruoming Pang. Conformer: Convolution-augmented transformer for speech recognition. In *Proceedings of the Conference of the International Speech Communication Association*, 2020.
- [17] Chun-Hua Guo and Nicholas J. Higham. A Schur–Newton method for the matrix p th root and its inverse. *SIAM Journal on Matrix Analysis and Applications*, 28(3):788–804, 2006.
- [18] Vineet Gupta, Tomer Koren, and Yoram Singer. Shampoo: Preconditioned stochastic tensor optimization. In *Proceedings of the International Conference on Machine Learning*, 2018.
- [19] N. Halko, P. G. Martinsson, and J. A. Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Review*, 53(2):217–288, 2011.
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, June 2016.
- [21] Roger A. Horn and Charles R. Johnson. *Matrix Analysis*. Cambridge university press, 2012.
- [22] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of the International Conference on Learning Representations*, 2015.
- [23] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [24] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- [25] Seung Hoon Lee, Seunghyun Lee, and Byung Cheol Song. Vision transformer for small-size datasets. *arXiv preprint arXiv:2112.13492*, 2021.
- [26] Bingrui Li, Jianfei Chen, and Jun Zhu. Memory efficient optimizers with 4-bit states. *Advances in Neural Information Processing Systems*, 2023.
- [27] Hong Liu, Zhiyuan Li, David Leo Wright Hall, Percy Liang, and Tengyu Ma. Sophia: A scalable stochastic second-order optimizer for language model pre-training. In *Proceedings of the International Conference on Learning Representations*, 2024.
- [28] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, October 2021.
- [29] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *Proceedings of the International Conference on Learning Representations*, 2019.
- [30] James Martens and Roger Grosse. Optimizing neural networks with Kronecker-factored approximate curvature. In *Proceedings of the International Conference on Machine Learning*, 2015.
- [31] Ning Qian. On the momentum term in gradient descent learning algorithms. *Neural Networks*, 12(1):145–151, 1999.
- [32] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, and others. Language models are unsupervised multitask learners. *OpenAI blog*, 2019.
- [33] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020.
- [34] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.

- [35] Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost. In *Proceedings of the International Conference on Machine Learning*, 2018.
- [36] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *Proceedings of the International Conference on Learning Representations*, 2015.
- [37] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, and others. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [38] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 2017.
- [39] Zhewei Yao, Amir Gholami, Sheng Shen, Mustafa Mustafa, Kurt Keutzer, and Michael W. Mahoney. AdaHessian: an adaptive second order optimizer for machine learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021.
- [40] Jui-Nan Yen, Sai Surya Duvvuri, Inderjit S. Dhillon, and Cho-Jui Hsieh. Block low-rank preconditioner with shared basis for stochastic optimization. *Advances in Neural Information Processing Systems*, 2023.
- [41] Hongwei Yong, Ying Sun, and Lei Zhang. A general regret bound of preconditioned gradient method for DNN training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, June 2023.
- [42] Lin Zhang, Shaohuai Shi, and Bo Li. Eva: Practical second-order optimization with Kronecker-vectorized approximation. In *Proceedings of the International Conference on Learning Representations*, 2023.
- [43] Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong Tian. GaLore: Memory-efficient LLM training by gradient low-rank projection. In *Proceedings of the International Conference on Machine Learning*, 2024.
- [44] Pan Zhou, Xingyu Xie, and Shuicheng Yan. Win: Weight-decay-integrated Nesterov acceleration for adaptive gradient algorithms. In *Proceedings of the International Conference on Learning Representations*, 2023.

A Implementation Details of Shampoo, CASPR, K-FAC and AdaBK

The implementation of 32-bit Shampoo used in our experiments is described in Algorithm 4. Our Pytorch implementation of Shampoo is partially based on the code provided by [2]. We implement CASPR by replacing $\hat{G}_t = \hat{L}_t G_t \hat{R}_t$ with $J_t = \hat{L}_t G_t + G_t \hat{R}_t$; $\hat{G}_t = \hat{L}_t J_t + J_t \hat{R}_t$ in line 12 of Algorithm 4 and line 14 of Algorithm 3. We summarize the implementation of 32-bit K-FAC/AdaBK in Algorithm 5, where X_t is the input feature and Y_t is the output feature gradient. Both power iteration [5] and Schur-Newton iteration [17] are run for 10 iterations. Our implementation of 4-bit K-FAC/AdaBK is similar to 4-bit Shampoo (i.e., compressing L_t , R_t , \hat{L}_t , and \hat{R}_t).

Algorithm 4 Practical 32-bit Shampoo

Input: initial parameter $W_0 \in \mathbb{R}^{m \times n}$, left preconditioner $L_0 = \epsilon I_m$, right preconditioner $R_0 = \epsilon I_n$, inverse root of left preconditioner $\hat{L}_0 = I_m$, inverse root of right preconditioner $\hat{R}_0 = I_n$, total number of steps T , interval of updating preconditioners T_1 , interval of updating inverse roots of preconditioners T_2 , exponential decay rate for preconditioners $\beta \in (0, 1)$, first-order optimizer \mathcal{F} , first-order optimizer state $s_0 = 0$.

Output: final parameter W_T .

```

1: for  $t = 1, 2, \dots, T$  do
2:   Receive loss function  $\mathcal{L}_t : \mathbb{R}^{m \times n} \mapsto \mathbb{R}$  and compute gradient  $G_t = \nabla \mathcal{L}_t(W_t)$ 
3:   if  $t \% T_1 \equiv 0$  then
4:      $L_t = \beta L_{t-1} + (1 - \beta) G_t G_t^\top$ ;  $R_t = \beta R_{t-1} + (1 - \beta) G_t^\top G_t$ 
5:   else
6:      $L_t = L_{t-1}$ ;  $R_t = R_{t-1}$ 
7:   if  $t \% T_2 \equiv 0$  then
8:     Compute maximum eigenvalues  $\lambda_{\max}^L$  and  $\lambda_{\max}^R$  of  $L_t$  and  $R_t$  by power iteration
9:     Compute  $\hat{L}_t = (L_t + \lambda_{\max}^L \epsilon I_m)^{-1/4}$  and  $\hat{R}_t = (R_t + \lambda_{\max}^R \epsilon I_n)^{-1/4}$  by Schur-Newton iteration
10:  else
11:     $\hat{L}_t = \hat{L}_{t-1}$ ;  $\hat{R}_t = \hat{R}_{t-1}$ 
12:     $\hat{G}_t = \hat{L}_t G_t \hat{R}_t$ ;  $\tilde{G}_t = \hat{G}_t (\|G_t\|_F / \|\hat{G}_t\|_F)$ 
13:     $W_t, s_t = \mathcal{F}(W_{t-1}, s_{t-1}, \tilde{G}_t)$ 

```

Algorithm 5 Practical 32-bit K-FAC/AdaBK

Input: initial parameter $W_0 \in \mathbb{R}^{m \times n}$, left preconditioner $L_0 = 0$, right preconditioner $R_0 = 0$, inverse root of left preconditioner $\hat{L}_0 = I_m$, inverse root of right preconditioner $\hat{R}_0 = I_n$, total number of steps T , interval of updating preconditioners T_1 , interval of updating inverse roots of preconditioners T_2 , ϵ , exponential decay rate for preconditioners $\beta \in (0, 1)$, $\alpha = 1$ for K-FAC / $\alpha = 2$ for AdaBK, first-order optimizer \mathcal{F} , first-order optimizer state $s_0 = 0$.

Output: final parameter W_T .

```

1: for  $t = 1, 2, \dots, T$  do
2:   Receive loss function  $\mathcal{L}_t : \mathbb{R}^{m \times n} \mapsto \mathbb{R}$  and compute gradient  $G_t = \nabla \mathcal{L}_t(W_t)$ 
3:   Receive  $X_t$  by forward propagation and  $Y_t$  by backward propagation
4:   if  $t \% T_1 \equiv 0$  then
5:      $L_t = \beta L_{t-1} + (1 - \beta) Y_t Y_t^\top$ ;  $R_t = \beta R_{t-1} + (1 - \beta) X_t X_t^\top$ 
6:   else
7:      $L_t = L_{t-1}$ ;  $R_t = R_{t-1}$ 
8:   if  $t \% T_2 \equiv 0$  then
9:     Compute maximum eigenvalues  $\lambda_{\max}^L$  and  $\lambda_{\max}^R$  of  $L_t$  and  $R_t$  by power iteration
10:    Compute  $\hat{L}_t = (L_t + \lambda_{\max}^L \epsilon I_m)^{-1/\alpha}$  and  $\hat{R}_t = (R_t + \lambda_{\max}^R \epsilon I_n)^{-1/\alpha}$  by Schur-Newton iteration
11:  else
12:     $\hat{L}_t = \hat{L}_{t-1}$ ;  $\hat{R}_t = \hat{R}_{t-1}$ 
13:     $\hat{G}_t = \hat{L}_t G_t \hat{R}_t$ ;  $\tilde{G}_t = \hat{G}_t (\|G_t\|_F / \|\hat{G}_t\|_F)$ 
14:     $W_t, s_t = \mathcal{F}(W_{t-1}, s_{t-1}, \tilde{G}_t)$ 

```

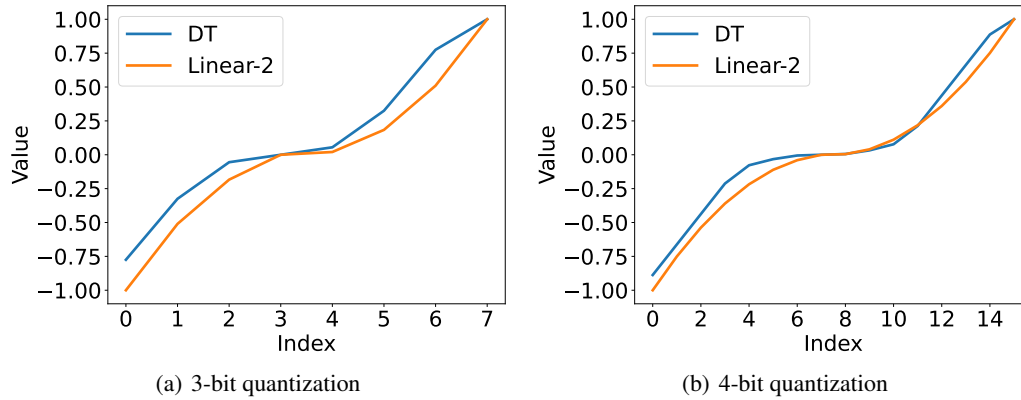


Figure 5: Visualization of DT quantization and Linear-2 quantization at b -bit ($b = 3, 4$) precision.

B Randomized SVD Method

Given an initial matrix $P_0 \in \mathbb{R}^{n \times n}$, randomized SVD method computes the eigenvector matrix of a PD matrix $A \in \mathbb{R}^{n \times n}$ by iterating

$$P_t = QR(AP_{t-1}), \quad (4)$$

where $QR(X)$ denotes the QR decomposition of matrix X , returning an orthogonal matrix. Since we can initialize P_0 with the previous result (e.g., V in Algorithm 1), only a few iterations are enough to obtain an accurate estimation in practice. In our experiments, we iterate (4) once for Shampoo/CASPR, and iterate (4) twice for K-FAC/AdaBK.

C Quantization Mappings

We present the constructions of different quantization mappings in b -bit quantizers (\mathcal{R} in \mathcal{Q}). See Figure 5 for the illustration of them. Note that $\mathbb{T}_b = \{0, 1, \dots, 2^b - 1\}$.

Dynamic tree (DT) quantization for b -bit quantization maps \mathbb{T}_b onto $\{0, 1\} \cup G$, where G is a set of numbers with the following properties: the number in G looks like $\pm q_k \times 10^{-E}$, where a) $b = 2 + E + F$, where E, F are integers; b) $q_k = (p_k + p_{k+1})/2$, where $k \in \{0, \dots, 2^F - 1\}$; c) $p_j = 0.9j/2^F + 0.1$, where $j \in \{0, \dots, 2^F\}$. For 4-bit quantization, DT quantization maps \mathbb{T}_4 onto $\{-0.8875, -0.6625, -0.4375, -0.2125, -0.0775, -0.0325, -0.0055, 0.0000, 0.0055, 0.0325, 0.0775, 0.2125, 0.4375, 0.6625, 0.8875, 1.0000\}$. For 3-bit quantization, DT quantization maps \mathbb{T}_3 onto $\{-0.7750, -0.3250, -0.0550, 0.0000, 0.0550, 0.3250, 0.7750, 1.0000\}$.

For 4-bit quantization, linear square (Linear-2) quantization maps \mathbb{T}_4 onto $\{-1.0000, -0.7511, -0.5378, -0.3600, -0.2178, -0.1111, -0.0400, 0.0000, 0.0044, 0.0400, 0.1111, 0.2178, 0.3600, 0.5378, 0.7511, 1.0000\}$. For 3-bit quantization, Linear-2 quantization maps \mathbb{T}_3 onto $\{-1.0000, -0.5102, -0.1837, 0.0000, 0.0204, 0.1837, 0.5102, 1.0000\}$.

D Quantization Error Analyses

We present more quantization error analyses of the preconditioners. Recall that we define two kinds of quantization errors in mapping f of transformation g at $A \in \mathbb{R}^{m \times n}$ (in short errors in $f(A)$ of g) in Subsection 3.1. Here we extend them as follows: define the normwise relative error (NRE) in f of (g_1, g_2) at A as

$$\text{NRE} = \frac{\|f(A) - g_2 \circ f \circ g_1(A)\|_F}{\|f(A)\|_F},$$

and the angle error (AE) in f of (g_1, g_2) at A as

$$\text{AE} = \arccos \left(\frac{\langle f(A), g_2 \circ f \circ g_1(A) \rangle}{\|f(A)\|_F \|g_2 \circ f \circ g_1(A)\|_F} \right).$$

D.1 Static Analysis

Table 5 is an extension of Table 1 for Bit=4. Since the diagonal elements of $\mathbf{A}^{-1/4}$ are usually much larger than its non-diagonal elements where \mathbf{A} is a PD matrix, we further consider the quantization errors in $f(\mathbf{A}) = \mathbf{A}^{-1/4} - \text{Diag}(\text{diag}(\mathbf{A}^{-1/4}))$ at 4-bit precision as shown in Table 6. Table 7 shows the quantization errors at 8-bit precision.

A large condition number of a PD matrix \mathbf{A} is indispensable for the superiority of quantizing \mathbf{U} over quantizing \mathbf{A} , where \mathbf{U} is the eigenvector matrix of \mathbf{A} . We consider contracting the singular value distribution of $\mathbf{A} = \mathbf{A}_1$ with SVD $\mathbf{U}\text{Diag}(\boldsymbol{\lambda})\mathbf{U}^\top$ used in Table 5 by mapping each singular value λ of \mathbf{A} to $h(\lambda) = \tau(\lambda - \lambda_{\min}^{\mathbf{A}}) + \lambda_{\min}^{\mathbf{A}}$, where $\lambda_{\min}^{\mathbf{A}}$ is the minimum singular value of \mathbf{A} and $\tau > 0$ is the contraction coefficient. Figure 6 shows 4-bit quantization errors in $\mathbf{A}^{-1/4}$ or $\mathbf{A}^{-1/4} - \text{Diag}(\text{diag}(\mathbf{A}^{-1/4}))$ of quantizing \mathbf{U} or \mathbf{A} at $\mathbf{A} = \mathbf{U}\text{Diag}(h(\boldsymbol{\lambda}))\mathbf{U}^\top$.

Table 5: Quantization errors in $f(\mathbf{A}) = \mathbf{A}^{-1/4}$ of different 4-bit quantization schemes at a PD matrix \mathbf{A} . We employ block-wise normalization with a block size of 64. \mathbf{U} is the eigenvector matrix of \mathbf{A} and $\mathbf{B} = (g_1(\mathbf{A}))^{-1/4}$. QM = quantized matrices and OR = orthogonal rectification.

Real-world $\mathbf{A} = \mathbf{A}_1$					Synthetic $\mathbf{A} = \mathbf{A}_2$				
Mapping \mathcal{R}	QM	OR	NRE ↓	AE (°) ↓	Mapping \mathcal{R}	QM	OR	NRE ↓	AE (°) ↓
DT	\mathbf{A}	✗	0.6241	17.319	DT	\mathbf{A}	✗	0.4615	17.189
	\mathbf{U}	✗	0.0709	4.0426		\mathbf{U}	✗	0.1224	7.0144
	\mathbf{U}	✓	0.0455	2.5615		\mathbf{U}	✓	0.0878	4.9960
	\mathbf{B}	✗	0.0398	2.2802		\mathbf{B}	✗	0.0853	4.8914
	(\mathbf{A}, \mathbf{B})	✗	0.6243	17.364		(\mathbf{A}, \mathbf{B})	✗	0.4649	17.650
	(\mathbf{U}, \mathbf{B})	✗	0.0811	4.6296		(\mathbf{U}, \mathbf{B})	✗	0.1485	8.5168
	(\mathbf{U}, \mathbf{B})	✓	0.0604	3.4230		(\mathbf{U}, \mathbf{B})	✓	0.1224	6.9817
Linear-2	\mathbf{A}	✗	0.6243	17.293	Linear-2	\mathbf{A}	✗	0.4465	15.338
	\mathbf{U}	✗	0.0543	3.1066		\mathbf{U}	✗	0.0942	5.3998
	\mathbf{U}	✓	0.0343	1.9456		\mathbf{U}	✓	0.0669	3.8166
	\mathbf{B}	✗	0.0315	1.8050		\mathbf{B}	✗	0.0661	3.7887
	(\mathbf{A}, \mathbf{B})	✗	0.6243	17.301		(\mathbf{A}, \mathbf{B})	✗	0.4483	15.654
	(\mathbf{U}, \mathbf{B})	✗	0.0626	3.5833		(\mathbf{U}, \mathbf{B})	✗	0.1150	6.5901
	(\mathbf{U}, \mathbf{B})	✓	0.0466	2.6494		(\mathbf{U}, \mathbf{B})	✓	0.0941	5.3716

Table 6: Quantization errors in $f(\mathbf{A}) = \mathbf{A}^{-1/4} - \text{Diag}(\text{diag}(\mathbf{A}^{-1/4}))$ of different 4-bit quantization schemes at a PD matrix \mathbf{A} . We employ block-wise normalization with a block size of 64. \mathbf{U} is the eigenvector matrix of \mathbf{A} and $\mathbf{B} = (g_1(\mathbf{A}))^{-1/4}$. QM = quantized matrices and OR = orthogonal rectification.

Real-world $\mathbf{A} = \mathbf{A}_1$					Synthetic $\mathbf{A} = \mathbf{A}_2$				
Mapping \mathcal{R}	QM	OR	NRE ↓	AE (°) ↓	Mapping \mathcal{R}	QM	OR	NRE ↓	AE (°) ↓
DT	\mathbf{A}	✗	0.9549	59.360	DT	\mathbf{A}	✗	0.6247	25.913
	\mathbf{U}	✗	0.2328	13.287		\mathbf{U}	✗	0.1994	11.444
	\mathbf{U}	✓	0.1480	8.4365		\mathbf{U}	✓	0.1427	8.1415
	\mathbf{B}	✗	0.1314	7.5513		\mathbf{B}	✗	0.1391	7.9813
	(\mathbf{A}, \mathbf{B})	✗	0.9561	59.825		(\mathbf{A}, \mathbf{B})	✗	0.6314	26.948
	(\mathbf{U}, \mathbf{B})	✗	0.2666	15.281		(\mathbf{U}, \mathbf{B})	✗	0.2420	13.911
	(\mathbf{U}, \mathbf{B})	✓	0.1977	11.322		(\mathbf{U}, \mathbf{B})	✓	0.1992	11.393
Linear-2	\mathbf{A}	✗	0.9547	58.336	Linear-2	\mathbf{A}	✗	0.6010	20.780
	\mathbf{U}	✗	0.1786	10.213		\mathbf{U}	✗	0.1534	8.8027
	\mathbf{U}	✓	0.1122	6.4096		\mathbf{U}	✓	0.1088	6.2176
	\mathbf{B}	✗	0.1041	5.9554		\mathbf{B}	✗	0.1078	6.1755
	(\mathbf{A}, \mathbf{B})	✗	0.9548	58.601		(\mathbf{A}, \mathbf{B})	✗	0.6047	21.666
	(\mathbf{U}, \mathbf{B})	✗	0.2063	11.778		(\mathbf{U}, \mathbf{B})	✗	0.1873	10.745
	(\mathbf{U}, \mathbf{B})	✓	0.1530	8.7337		(\mathbf{U}, \mathbf{B})	✓	0.1532	8.7534

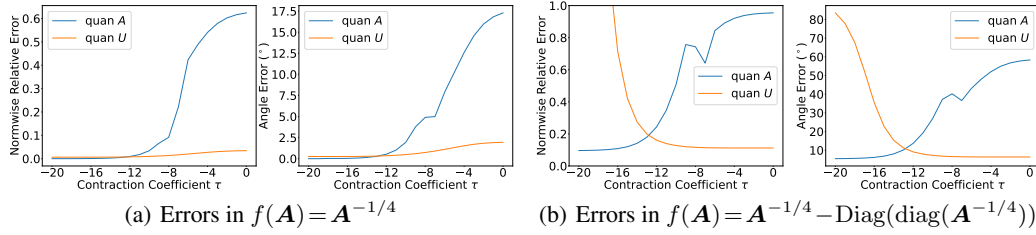


Figure 6: 4-bit quantization errors in $f(\mathbf{A})$ of quantizing \mathbf{U} or \mathbf{A} at $\mathbf{A} = \mathbf{U}\text{Diag}(h(\boldsymbol{\lambda}))\mathbf{U}^\top$. We use linear square quantization and orthogonal rectification. The condition number $\text{cond}(\mathbf{A}) = \lambda_{\max}^A / \lambda_{\min}^A$ is around 37235, where λ_{\max}^A and λ_{\min}^A are the maximum and minimum singular values of \mathbf{A} respectively. Contraction coefficients are shown on a \log_2 scale.

Table 7: Quantization errors in $f(\mathbf{A})$ of different 8-bit quantization schemes at a PD matrix \mathbf{A} , where $\mathbf{A} = \mathbf{A}_1$ is derived from the real world as described in Subsection 3.1. We employ block-wise normalization with a block size of 256. \mathbf{U} is the eigenvector matrix of \mathbf{A} and $\mathbf{B} = (g_1(\mathbf{A}))^{-1/4}$. QM = quantized matrices and OR = orthogonal rectification.

$f(\mathbf{A}) = \mathbf{A}^{-1/4}$					$f(\mathbf{A}) = \mathbf{A}^{-1/4} - \text{Diag}(\text{diag}(\mathbf{A}^{-1/4}))$				
Mapping \mathcal{R}	QM	OR	NRE \downarrow	AE ($^\circ$) \downarrow	Mapping \mathcal{R}	QM	OR	NRE \downarrow	AE ($^\circ$) \downarrow
DT	\mathbf{A}	\times	0.2192	8.3014	DT	\mathbf{A}	\times	0.5001	23.644
	\mathbf{U}	\times	0.0060	0.3421		\mathbf{U}	\times	0.0197	1.1273
	\mathbf{U}	\checkmark	0.0037	0.2140		\mathbf{U}	\checkmark	0.0123	0.7022
	\mathbf{B}	\times	0.0029	0.1655		\mathbf{B}	\times	0.0097	0.5553
	(\mathbf{A}, \mathbf{B})	\times	0.2193	8.3051		(\mathbf{A}, \mathbf{B})	\times	0.5003	23.649
	(\mathbf{U}, \mathbf{B})	\times	0.0067	0.3810		(\mathbf{U}, \mathbf{B})	\times	0.0219	1.2577
	(\mathbf{U}, \mathbf{B})	\checkmark	0.0047	0.2712		(\mathbf{U}, \mathbf{B})	\checkmark	0.0156	0.8955
Linear-2	\mathbf{A}	\times	0.2164	7.9751	Linear-2	\mathbf{A}	\times	0.4875	21.447
	\mathbf{U}	\times	0.0037	0.2121		\mathbf{U}	\times	0.0122	0.6994
	\mathbf{U}	\checkmark	0.0023	0.1312		\mathbf{U}	\checkmark	0.0076	0.4343
	\mathbf{B}	\times	0.0021	0.1203		\mathbf{B}	\times	0.0070	0.4035
	(\mathbf{A}, \mathbf{B})	\times	0.2164	7.9755		(\mathbf{A}, \mathbf{B})	\times	0.4875	21.448
	(\mathbf{U}, \mathbf{B})	\times	0.0043	0.2439		(\mathbf{U}, \mathbf{B})	\times	0.0141	0.8079
	(\mathbf{U}, \mathbf{B})	\checkmark	0.0031	0.1791		(\mathbf{U}, \mathbf{B})	\checkmark	0.0104	0.5935

D.2 Dynamic Analysis

We define the normwise relative error (NRE) and angle error (AE) of \mathbf{B} deviating from \mathbf{A} as

$$\text{NRE} = \frac{\|\mathbf{B} - \mathbf{A}\|_F}{\|\mathbf{A}\|_F}, \quad \text{AE} = \arccos \left(\frac{\langle \mathbf{A}, \mathbf{B} \rangle}{\|\mathbf{A}\|_F \|\mathbf{B}\|_F} \right).$$

Consider Shampoo using 4-bit preconditioners for parameter updates, but also recording 32-bit preconditioners at the same time. We extract the left preconditioners \mathbf{L}_4 and $\mathbf{L}_{32} \in \mathbb{R}^{1200 \times 1200}$ of a specific model parameter block $\mathbf{W} \in \mathbb{R}^{1200 \times 768}$ every 8000 steps in the Swin-Tiny training on CIFAR-100 with AdamW+Shampoo. Here \mathbf{L}_4 is a decompressed 4-bit preconditioner, and \mathbf{L}_{32} is a 32-bit preconditioner.

Figure 7 shows the quantization errors during training. For naive 4-bit Shampoo, $\mathbf{L}_{32}^{-1/4}$ and $\mathbf{L}_4^{-1/4}$ are computed by Schur-Newton iteration used in Algorithm 4 where $\epsilon = 10^{-4}$. For our 4-bit Shampoo, $\mathbf{L}_{32}^{-1/4}$ is computed by Schur-Newton iteration used in Algorithm 4 where $\epsilon = 10^{-4}$, and $\mathbf{L}_4^{-1/4}$ is computed by Algorithm 2 without quantization where $\epsilon = 10^{-4}$, $t_2 = 4$. We find that $\epsilon = 10^{-6}$ for Algorithm 2 used in our main experiments though is effective, yet it can cause a large numerical instability in the later stage of training (see Figure 8).

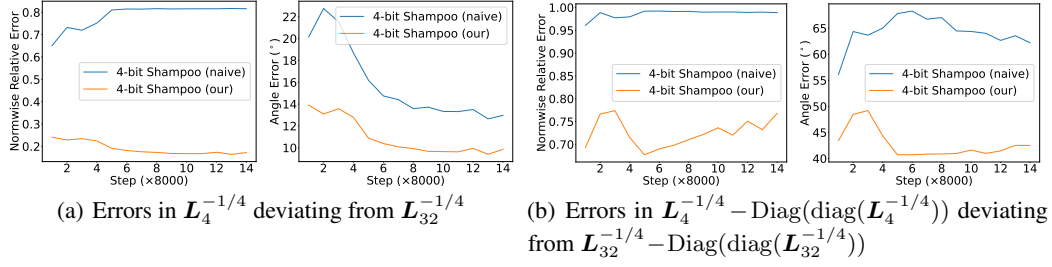


Figure 7: Quantization errors during Swin-Tiny training on the CIFAR-100 dataset. We use dampening term $\epsilon = 10^{-4}$ to compute $\mathbf{L}_4^{-1/4}$ and $\mathbf{L}_{32}^{-1/4}$.

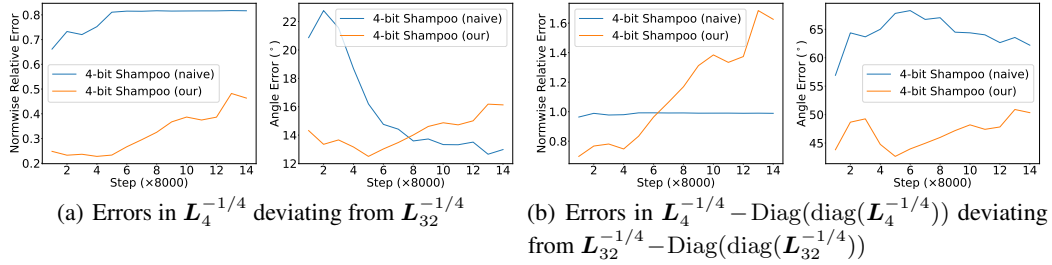


Figure 8: Quantization errors during Swin-Tiny training on the CIFAR-100 dataset. We use dampening term $\epsilon = 10^{-6}$ to compute $\mathbf{L}_4^{-1/4}$ and $\mathbf{L}_{32}^{-1/4}$.

E Convergence Analysis

More notations. Given a symmetric real matrix \mathbf{A} , $\mathbf{A} \succeq 0$ means that \mathbf{A} is positive semidefinite (PSD), and $\mathbf{A} \succ 0$ means that \mathbf{A} is positive definite (PD). Assume that symmetric matrices \mathbf{A} and \mathbf{B} are symmetric, the notations $\mathbf{A} \succeq \mathbf{B}$ and $\mathbf{A} \succ \mathbf{B}$ mean that $\mathbf{A} - \mathbf{B} \succeq 0$ and $\mathbf{A} - \mathbf{B} \succ 0$ respectively. Let \mathbf{A} be a PSD matrix and $s \in \mathbb{R}$, we define $\mathbf{A}^s = \mathbf{U}\mathbf{\Lambda}^s\mathbf{U}^\top$, where $\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$ is the Singular Value Decomposition (SVD) of \mathbf{A} . The Mahalanobis norm of a vector \mathbf{x} induced by a PD matrix \mathbf{A} is $\|\mathbf{x}\|_{\mathbf{A}} = \sqrt{\mathbf{x}^\top \mathbf{A} \mathbf{x}}$. The dual norm of $\|\cdot\|_{\mathbf{A}}$ is denoted by $\|\cdot\|_{\mathbf{A}}^*$, where $\|\mathbf{x}\|_{\mathbf{A}}^* = \sqrt{\mathbf{x}^\top \mathbf{A}^{-1} \mathbf{x}}$. The spectral norm of matrix \mathbf{A} is $\|\mathbf{A}\|_2 = \sup_{\mathbf{x} \neq 0} \{\|\mathbf{A}\mathbf{x}\|_2 / \|\mathbf{x}\|_2\}$. $\mathbf{A} \otimes \mathbf{B}$ means the (right) Kronecker product of matrices \mathbf{A} and \mathbf{B} . $\text{vec}(\mathbf{A})$ means the vectorization (stacking the rows) of \mathbf{A} .

Algorithm 6 Perturbed Shampoo in the matrix case

Input: $\mathbf{W}_0 \in \mathbb{R}^{m \times n}$, $\mathbf{L}_0 = \mathbf{0}_{m \times m}$, $\mathbf{R}_0 = \mathbf{0}_{n \times n}$, $\rho_0 = 0$, $\mu_0 = 0$.

1: **for** $t = 1, \dots, T$ **do**

2: Receive loss function: $f_t : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$

3: Compute gradient: $\mathbf{G}_t = \nabla f_t(\mathbf{W}_t)$

4: Update preconditioners: $\mathbf{J}_t = \mathbf{L}_{t-1} + \mathbf{G}_t \mathbf{G}_t^\top$; $\mathbf{K}_t = \mathbf{R}_{t-1} + \mathbf{G}_t^\top \mathbf{G}_t$

5: Perturb preconditioners: $\mathbf{L}_t = g(\mathbf{J}_t)$; $\mathbf{R}_t = g(\mathbf{K}_t)$

6: Accumulate errors: $\rho_t = \rho_{t-1} + \|\mathbf{J}_t - \mathbf{L}_t\|_2$; $\mu_t = \mu_{t-1} + \|\mathbf{K}_t - \mathbf{R}_t\|_2$

7: Update parameters: $\mathbf{W}_{t+1} = \mathbf{W}_t - \eta((\epsilon + \rho_t)\mathbf{I}_m + \mathbf{L}_t)^{-1/4} \mathbf{G}_t ((\epsilon + \mu_t)\mathbf{I}_n + \mathbf{R}_t)^{-1/4}$

We consider quantization as a perturbation and present the perturbed Shampoo in Algorithm 6 for convergence analysis. The regret bound of the perturbed Shampoo can be found in Theorem 1. Complete proofs can be found in Appendix F. We first introduce some basic technical tools, and the details of them are in [18, 21].

Lemma 3. Let $\mathbf{A}, \mathbf{A}', \mathbf{B}, \mathbf{B}'$ be matrices of appropriate dimensions, and \mathbf{u}, \mathbf{v} be two column vectors. The following properties hold:

(1) $(\mathbf{A} \otimes \mathbf{B})(\mathbf{A}' \otimes \mathbf{B}') = (\mathbf{A}\mathbf{A}') \otimes (\mathbf{B}\mathbf{B}')$;

- (2) $(A \otimes B)^\top = (A^\top \otimes B^\top)$;
(3) If $A, B \succeq 0$ and $s \in \mathbb{R}$, then $(A \otimes B)^s = (A^s \otimes B^s)$;
(4) If $A \succeq A'$ and $B \succeq B'$, then $A \otimes B \succeq A' \otimes B'$;
(5) $\text{tr}(AB) = \text{tr}(A)\text{tr}(B)$;
(6) $\text{vec}(uv^\top) = u \otimes v$.

Lemma 4. Let $G \in \mathbb{R}^{m \times n}$, $L \in \mathbb{R}^{m \times m}$, $R \in \mathbb{R}^{n \times n}$, then it holds that

$$(L \otimes R^\top) \text{vec}(G) = \text{vec}(LGR).$$

Lemma 5. Assume that $0 \preceq X_i \preceq Y_i$ for $i = 1, \dots, n$. Assume further that all X_i commute with each other and all Y_i commute with each other. Let $\alpha_1, \dots, \alpha_n \geq 0$ such that $\sum_{i=1}^n \alpha_i = 1$, then

$$X_1^{\alpha_1} \dots X_n^{\alpha_n} \preceq Y_1^{\alpha_1} \dots Y_n^{\alpha_n}.$$

Lemma 6. Let $0 \leq \alpha \leq 1$ and $0 \preceq X \preceq Y$, then $X^\alpha \preceq Y^\alpha$.

Lemma 7. Let $A \succ 0$ and $B \succ 0$, then it holds that $A \succeq B$ if and only if $B^{-1} \succeq A^{-1}$.

Lemma 8 (von Neumann). Let $A, B \in \mathbb{R}^{m \times n}$ and $q = \min\{m, n\}$. Let $\sigma_1(A) \geq \dots \geq \sigma_q(A)$ and $\sigma_1(B) \geq \dots \geq \sigma_q(B)$ denote the non-increasingly ordered singular values of A and B , respectively. Then

$$\langle A, B \rangle \leq \sum_{i=1}^q \sigma_i(A) \sigma_i(B).$$

Lemma 9. Assume that function f_t is continuously differentiable and convex on \mathbb{R}^d , and matrix $H_t \succ 0$ for $t = 1, \dots, T$. Given $w_0 \in \mathbb{R}^d$, $\eta > 0$, define $w_{t+1} = w_t - \eta H_t^{-1} g_t$, where $g_t = \nabla f_t(w_t)$. Then for any $w^* \in \mathbb{R}^d$, we have

$$\sum_{t=1}^T f_t(w_t) - \sum_{t=1}^T f_t(w^*) \leq \frac{1}{2\eta} \sum_{t=1}^T (\|w_t - w^*\|_{H_t}^2 - \|w_{t+1} - w^*\|_{H_t}^2) + \frac{\eta}{2} \sum_{t=1}^T (\|g_t\|_{H_t}^*)^2.$$

Lemma 10. Let g_1, \dots, g_T be a sequence of vectors. For $\rho > 0$, define $\widehat{H}_t = (\rho I + \sum_{s=1}^t g_s g_s^\top)^{1/2}$. Then we have

$$\sum_{t=1}^T (\|g_t\|_{\widehat{H}_t}^*)^2 \leq 2\text{tr}(\widehat{H}_T).$$

Lemma 11. Assume that $G_1, \dots, G_T \in \mathbb{R}^{m \times n}$ are matrices of rank at most r . Let s for $t = 1, \dots, T$. Then for any $\epsilon \geq 0$,

$$\epsilon I_{mn} + \frac{1}{r} \sum_{t=1}^T g_t g_t^\top \preceq (\epsilon I_m + \sum_{t=1}^T G_t G_t^\top)^{1/2} \otimes (\epsilon I_n + \sum_{t=1}^T G_t^\top G_t)^{1/2}.$$

The key to the convergence proof of Algorithm 6 is forming a PD matrix sequence $\{H_i\}_{i=1}^T$, which satisfies $0 \prec H_1 \preceq \dots \preceq H_T$. To achieve it, we give the following lemma extended from Lemma 2 in the Appendix of [40].

Lemma 12. Let $\{X_t\}_{t=1}^T$ be a sequence of symmetric matrices, and $A_t = \sum_{s=1}^t X_s$, where $t = 1, \dots, T$. Suppose we have two sequences of symmetric matrices $\{Y_t\}_{t=1}^T, \{Z_t\}_{t=0}^T$, and a sequence real numbers $\{\rho_t\}_{t=0}^T$ satisfying

$$Y_t = Z_{t-1} + X_t, \quad \rho_t = \rho_{t-1} + \|Y_t - Z_t\|_2, \quad Z_0 = 0, \rho_0 = 0.$$

Define $B_t = \rho_t I + Z_t$, where I denotes the identity matrix. Then for $t = 1, \dots, T$, we have

$$B_t \succeq B_{t-1} + X_t, \quad A_t \preceq B_t \preceq 2\rho_t I + A_t.$$

Theorem 1. Assume that the gradients $G_1, \dots, G_T \in \mathbb{R}^{m \times n}$ are matrices of rank at most r . Then for any $W^* \in \mathbb{R}^{m \times n}$ and $\epsilon > 0$, if $\eta = D/\sqrt{2r}$, the regret of Algorithm 6 is bounded as follows,

$$\sum_{t=1}^T f_t(W_t) - \sum_{t=1}^T f_t(W^*) \leq \sqrt{2r}D[2^{1/4}m\rho_T^{1/4} + \text{tr}(\tilde{L}_T^{1/4})][2^{1/4}n\mu_T^{1/4} + \text{tr}(\tilde{R}_T^{1/4})],$$

where $D = \max_{t \in [T]} \|W_t - W^*\|_F$, $\tilde{L}_t = \epsilon I_m + \sum_{s=1}^t G_s G_s^\top$, and $\tilde{R}_t = \epsilon I_n + \sum_{s=1}^t G_s^\top G_s$.

Though we get a convergence guarantee of Algorithm 6, the upper bound given by Theorem 1 is very slack, since $2^{1/4}m\rho_T^{1/4}$ is about the same as $\text{tr}(\tilde{L}_T^{1/4})$ for 4-bit quantization schemes in practice.

F Proofs

Lemma 1. Let \mathbf{A} be a PD matrix whose SVD is $\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$, where $\mathbf{U}=[\mathbf{u}_i]$ is an orthogonal matrix and $\mathbf{\Lambda}=\text{diag}([\lambda_i]^\top)$ is a diagonal matrix. Given a perturbation $\Delta\mathbf{U}=[\Delta\mathbf{u}_i]$ and $s \in \mathbb{R}$, we define $\mathbf{B}:=(\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top)^s$ and $\Delta\mathbf{B}:=((\mathbf{U}+\Delta\mathbf{U})\mathbf{\Lambda}(\mathbf{U}+\Delta\mathbf{U})^\top)^s-\mathbf{B}$.

(1) If $\mathbf{U}+\Delta\mathbf{U}$ is orthogonal and there exists $\alpha \in \mathbb{R}$ such that $\|\Delta\mathbf{u}_i\|_2 \leq \alpha$, then

$$\frac{\|\Delta\mathbf{B}\|_F}{\|\mathbf{B}\|_F} \leq 2\alpha.$$

(2) If $\mathbf{U}+\Delta\mathbf{U}$ is orthogonal and there exists $\beta \in \mathbb{R}$ such that $\langle \mathbf{u}_i, \mathbf{u}_i + \Delta\mathbf{u}_i \rangle \geq 1 - \beta \geq 0$, then

$$\frac{\langle \mathbf{B}, \mathbf{B} + \Delta\mathbf{B} \rangle}{\|\mathbf{B}\|_F \|\mathbf{B} + \Delta\mathbf{B}\|_F} \geq (1 - \beta)^2.$$

Proof. (1) Since \mathbf{U} and $\mathbf{U} + \Delta\mathbf{U}$ are orthogonal, we have

$$\mathbf{B} = \mathbf{U}\mathbf{\Lambda}^s\mathbf{U}^\top, \quad \mathbf{B} + \Delta\mathbf{B} = (\mathbf{U} + \Delta\mathbf{U})\mathbf{\Lambda}^s(\mathbf{U} + \Delta\mathbf{U})^\top,$$

by definition. This leads to

$$\Delta\mathbf{B} = \mathbf{U}\mathbf{\Lambda}^s\Delta\mathbf{U}^\top + \Delta\mathbf{U}\mathbf{\Lambda}^s(\mathbf{U} + \Delta\mathbf{U})^\top.$$

The Frobenius norm satisfies the triangle inequality and is orthogonality invariant. Hence,

$$\begin{aligned} \|\Delta\mathbf{B}\|_F &= \|\mathbf{U}\mathbf{\Lambda}^s\Delta\mathbf{U}^\top + \Delta\mathbf{U}\mathbf{\Lambda}^s(\mathbf{U} + \Delta\mathbf{U})^\top\|_F \\ &\leq \|\mathbf{U}\mathbf{\Lambda}^s\Delta\mathbf{U}^\top\|_F + \|\Delta\mathbf{U}\mathbf{\Lambda}^s(\mathbf{U} + \Delta\mathbf{U})^\top\|_F \\ &= \|\mathbf{\Lambda}^s\Delta\mathbf{U}^\top\|_F + \|\Delta\mathbf{U}\mathbf{\Lambda}^s\|_F = 2\|\Delta\mathbf{U}\mathbf{\Lambda}^s\|_F \\ &= 2\sqrt{\sum_i \lambda_i^{2s} \|\Delta\mathbf{u}_i\|_2^2} = 2\sqrt{\sum_i \lambda_i^{2s} \|\Delta\mathbf{u}_i\|_2^2} \\ &\leq 2\sqrt{\sum_i \lambda_i^{2s} \alpha^2} = 2\alpha\sqrt{\sum_i \lambda_i^{2s}} = 2\alpha\|\mathbf{\Lambda}^s\|_F \\ &= 2\alpha\|\mathbf{B}\|_F. \end{aligned}$$

(2) Similar to (1), we have

$$\Delta\mathbf{B} = \mathbf{U}\mathbf{\Lambda}^s\Delta\mathbf{U}^\top + \Delta\mathbf{U}\mathbf{\Lambda}^s\mathbf{U}^\top + \Delta\mathbf{U}\mathbf{\Lambda}^s\Delta\mathbf{U}^\top.$$

From $\langle \mathbf{u}_i, \mathbf{u}_i + \Delta\mathbf{u}_i \rangle \geq 1 - \beta \geq 0$, we get $0 \geq \langle \mathbf{u}_i, \Delta\mathbf{u}_i \rangle \geq -\beta \geq -1$ because

$$1 = \|\mathbf{u}_i\|_2 \|\mathbf{u}_i + \Delta\mathbf{u}_i\|_2 \geq \langle \mathbf{u}_i, \mathbf{u}_i + \Delta\mathbf{u}_i \rangle = 1 + \langle \mathbf{u}_i, \Delta\mathbf{u}_i \rangle \geq 1 - \beta \geq 0,$$

holds due to the orthogonality of \mathbf{U} and $\mathbf{U} + \Delta\mathbf{U}$. Hence,

$$\begin{aligned} \langle \mathbf{B}, \Delta\mathbf{B} \rangle &= \text{tr}(2\mathbf{U}\mathbf{\Lambda}^{2s}\Delta\mathbf{U}^\top + \mathbf{U}\mathbf{\Lambda}^s\mathbf{U}^\top\Delta\mathbf{U}\mathbf{\Lambda}^s\Delta\mathbf{U}^\top) \\ &= \text{tr}\left(\sum_i 2\lambda_i^{2s} \mathbf{u}_i \Delta\mathbf{u}_i^\top\right) + \text{tr}\left[\left(\sum_i \lambda_i^s \mathbf{u}_i \mathbf{u}_i^\top\right) \left(\sum_j \lambda_j^s \Delta\mathbf{u}_j \Delta\mathbf{u}_j^\top\right)\right] \\ &= \left(\sum_i 2\lambda_i^{2s} \langle \mathbf{u}_i, \Delta\mathbf{u}_i \rangle\right) + \left(\sum_{ij} \lambda_i^s \lambda_j^s \langle \mathbf{u}_i, \Delta\mathbf{u}_j \rangle^2\right) \\ &\geq \left(\sum_i 2\lambda_i^{2s} \langle \mathbf{u}_i, \Delta\mathbf{u}_i \rangle\right) + \left(\sum_i \lambda_i^{2s} \langle \mathbf{u}_i, \Delta\mathbf{u}_i \rangle^2\right) \\ &= \sum_i \lambda_i^{2s} [(1 + \langle \mathbf{u}_i, \Delta\mathbf{u}_i \rangle)^2 - 1] \\ &\geq \sum_i \lambda_i^{2s} [(1 - \beta)^2 - 1] = [(1 - \beta)^2 - 1] \|\mathbf{\Lambda}^s\|_F^2 \\ &= [(1 - \beta)^2 - 1] \|\mathbf{B}\|_F^2 = [(1 - \beta)^2 - 1] \langle \mathbf{B}, \mathbf{B} \rangle. \end{aligned}$$

Therefore, we have

$$\frac{\langle \mathbf{B}, \mathbf{B} + \Delta\mathbf{B} \rangle}{\|\mathbf{B}\|_F \|\mathbf{B} + \Delta\mathbf{B}\|_F} = \frac{\langle \mathbf{B}, \mathbf{B} + \Delta\mathbf{B} \rangle}{\langle \mathbf{B}, \mathbf{B} \rangle} = 1 + \frac{\langle \mathbf{B}, \Delta\mathbf{B} \rangle}{\langle \mathbf{B}, \mathbf{B} \rangle} \geq (1 - \beta)^2.$$

The proof is completed. \square

Lemma 2. Let \mathbf{A} be a PD matrix of order $m+n$ whose SVD is $\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$, where $m, n \in \mathbb{N}_+$, $n = lm$, $\mathbf{U} = [\mathbf{u}_i]$ is an orthogonal matrix and $\mathbf{\Lambda} = \text{diag}([\lambda_i]^\top)$ is a diagonal matrix. Assume that $\mathbf{\Lambda} = \text{diag}([c\lambda\mathbf{1}_{m \times 1}^\top, \lambda\mathbf{1}_{n \times 1}^\top]^\top)$, $c \geq 1$, and $\lambda > 0$. Given a perturbation $\Delta\mathbf{\Lambda} = \text{diag}([\mathbf{0}_{m \times 1}^\top, \Delta\lambda_{n \times 1}^\top]^\top)$ and $s \in \mathbb{R}$, we define $\mathbf{B} := (\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top)^s$ and $\Delta\mathbf{B} := (\mathbf{U}(\mathbf{\Lambda} + \Delta\mathbf{\Lambda})\mathbf{U}^\top)^s - \mathbf{B}$.

(1) If $\Delta\lambda_{n \times 1} = (k-1)\lambda\mathbf{1}_{n \times 1}$ where $k > 0$, then

$$\frac{\|\Delta\mathbf{B}\|_F}{\|\mathbf{B}\|_F} = \frac{\sqrt{l}|k^s - 1|}{\sqrt{c^{2s} + l}} = h_1(s, l).$$

Moreover, $h_1(s, l)$ decreases monotonically with s over $(-\infty, 0)$ and increases monotonically with l over $(0, +\infty)$.

(2) If $\Delta\lambda_{n \times 1} = (tc-1)\lambda\mathbf{1}_{n \times 1}$ where $t > 0$, then

$$\frac{\langle \mathbf{B}, \mathbf{B} + \Delta\mathbf{B} \rangle}{\|\mathbf{B}\|_F \|\mathbf{B} + \Delta\mathbf{B}\|_F} = \frac{lt^s + c^s}{\sqrt{(1 + lt^{2s})(l + c^{2s})}} = h_2(l).$$

Moreover, $h_2(l)$ decreases monotonically with l over $(0, (c/t)^s]$ and increases monotonically with l over $((c/t)^s, +\infty)$.

(3) If $\Delta\lambda_{n \times 1} = (tc-1)\lambda\mathbf{1}_{n \times 1}$ where $k = tc > 0$ and $l = (c/t)^s$, then

$$\frac{\|\Delta\mathbf{B}\|_F}{\|\mathbf{B}\|_F} = \frac{|k^s - 1|}{\sqrt{k^s + 1}}, \quad \frac{\langle \mathbf{B}, \mathbf{B} + \Delta\mathbf{B} \rangle}{\|\mathbf{B}\|_F \|\mathbf{B} + \Delta\mathbf{B}\|_F} = \frac{2}{\sqrt{2 + k^s + 1/k^s}}.$$

Proof. (1) Since \mathbf{U} is orthogonal, we have

$$\|\Delta\mathbf{B}\|_F = \|(\mathbf{\Lambda} + \Delta\mathbf{\Lambda})^s - \mathbf{\Lambda}^s\|_F = \sqrt{n}|k^s - 1|\lambda^s, \quad \|\mathbf{B}\|_F = \|\mathbf{\Lambda}^s\|_F = \sqrt{mc^{2s} + n\lambda^s}.$$

Hence,

$$\frac{\|\Delta\mathbf{B}\|_F}{\|\mathbf{B}\|_F} = \frac{\sqrt{n}|k^s - 1|}{\sqrt{mc^{2s} + n}} = \frac{\sqrt{l}|k^s - 1|}{\sqrt{c^{2s} + l}} = h_1(s, l) \geq 0.$$

It is easy to check that h_1 increases monotonically with l over $(0, +\infty)$. To prove h_1 decreases monotonically with s over $(-\infty, 0)$, define

$$g_1(s) = \frac{1}{l}(h_1(s, l))^2 = \frac{(k^s - 1)^2}{c^{2s} + l}.$$

Consider the derivative of g_1

$$\begin{aligned} g'_1(s) &= \frac{(c^{2s} + l)2(k^s - 1)k^s \ln k - (k^s - 1)^2 c^{2s} 2 \ln c}{(c^{2s} + l)^2} \\ &= \frac{2(k^s - 1)((c^{2s} + l)k^s \ln k - (k^s - 1)c^{2s} \ln c)}{(c^{2s} + l)^2}. \end{aligned}$$

If $s < 0$ and $k > 1$, then $k^s - 1 < 0$, $k^s \ln k > 0$ leading to $g'_1(s) < 0$ since $c \geq 0$; Similarly, if $s < 0$ and $0 < k \leq 1$, then $k^s - 1 \geq 0$, $k^s \ln k \leq 0$ leading to $g'_1(s) \leq 0$. Thus $g_1(s)$ is a monotonically decreasing function for $s < 0$, which implies that h_1 decreases monotonically with s over $(-\infty, 0)$.

(2) Similar to (1), we have

$$\|\mathbf{B}\|_F = \sqrt{mc^{2s} + n\lambda^s}, \quad \|\mathbf{B} + \Delta\mathbf{B}\|_F = \sqrt{nt^{2s} + mc^s\lambda^s}.$$

Besides,

$$\langle \mathbf{B}, \mathbf{B} + \Delta\mathbf{B} \rangle = \text{tr}(\mathbf{U}\mathbf{\Lambda}^s(\mathbf{\Lambda} + \Delta\mathbf{\Lambda})^s\mathbf{U}^\top) = \text{tr}(\mathbf{\Lambda}^s(\mathbf{\Lambda} + \Delta\mathbf{\Lambda})^s) = (mc^{2s} + nc^s t^s)\lambda^{2s}.$$

Hence, we get

$$\frac{\langle \mathbf{B}, \mathbf{B} + \Delta\mathbf{B} \rangle}{\|\mathbf{B}\|_F \|\mathbf{B} + \Delta\mathbf{B}\|_F} = \frac{nt^s + mc^s}{\sqrt{(m + nt^{2s})(n + mc^{2s})}} = \frac{lt^s + c^s}{\sqrt{(1 + lt^{2s})(l + c^{2s})}} = h_2(l) \geq 0.$$

To prove h_2 decreases monotonically with l over $(0, (c/t)^s]$ and increases monotonically with l over $((c/t)^s, +\infty)$, we define

$$g_2(l) = (h_2(l))^2 = \frac{(lt^s + c^s)^2}{(1 + lt^{2s})(l + c^{2s})},$$

whose monotonicity is equivalent to that of h_2 for $l > 0$. Consider the derivative of g_2

$$g_2'(l) = \left(\frac{t^{2s}l^2 + 2t^sc^sl + c^{2s}}{t^{2s}l^2 + l + t^{2s}c^{2s}l + c^{2s}} \right)' = \frac{(t^s - t^{2s}c^s)^2l^2 - (c^s - t^sc^{2s})^2}{(t^{2s}l^2 + l + t^{2s}c^{2s}l + c^{2s})^2}.$$

If $s = 0$ or $tc = 1$, then $g_2(l) \equiv 1$. If $s \neq 0$ and $tc \neq 1$, then $(t^s - t^{2s}c^s)^2 > 0$, $(c^s - t^sc^{2s})^2 > 0$. In this case, let $g_2'(l) = 0$, we get

$$t^{2s}(1 - t^sc^s)^2l^2 = c^{2s}(1 - t^sc^s)^2,$$

which implies that $l = (c/t)^s$. It is easy to see that g_2 decreases monotonically with l over $(0, (c/t)^s]$ and increases monotonically with l over $((c/t)^s, +\infty)$.

(3) According to (1)(2), we can easily get

$$\frac{\|\Delta B\|_F}{\|B\|_F} = \frac{|k^s - 1|}{\sqrt{k^s + 1}}, \quad \frac{\langle B, B + \Delta B \rangle}{\|B\|_F \|B + \Delta B\|_F} = \frac{2}{\sqrt{2 + k^s + 1/k^s}}.$$

The proof is completed. \square

Proposition 1. Let A be a PD matrix of order $m+n$ whose SVD is $U\Lambda U^T$, where $m, n \in \mathbb{N}_+$, $n = lm$, $U = [u_i]$ is an orthogonal matrix, $\Lambda = \text{diag}([c\lambda 1_{m \times 1}^T, \lambda 1_{n \times 1}^T]^T)$, $c \geq 1000$, and $\lambda > 0$. Given $\Delta U = [\Delta u_i]$, $\Delta \Lambda = \text{diag}([0_{m \times 1}^T, \Delta \lambda_{n \times 1}^T]^T)$, and $s \leq -0.25$, we define $B := (U\Lambda U^T)^s$, $B_1 := ((U + \Delta U)\Lambda(U + \Delta U)^T)^s$, and $B_2 := (U(\Lambda + \Delta \Lambda)U^T)^s$. If $U + \Delta U$ is orthogonal, $\|\Delta u_i\|_2 \leq 0.1$, $\langle u_i, \Delta u_i \rangle \geq -0.005$, $\Delta \lambda_{n \times 1} = (0.02c - 1)\lambda 1_{n \times 1}$, and $l = (c/0.02)^s$, then

$$2 \frac{\|B_1 - B\|_F}{\|B\|_F} \leq 0.4 \leq \frac{\|B_2 - B\|_F}{\|B\|_F}, \quad 6 \left(1 - \frac{\langle B, B_1 \rangle}{\|B\|_F \|B_1\|_F} \right) \leq 0.06 \leq \left(1 - \frac{\langle B, B_2 \rangle}{\|B\|_F \|B_2\|_F} \right).$$

Proof. According to Lemma 1, we have

$$\frac{\|B_1 - B\|_F}{\|B\|_F} \leq 0.2, \quad \frac{\langle B, B_1 \rangle}{\|B\|_F \|B_1\|_F} \geq (1 - 0.005)^2 \geq 0.99.$$

On the other hand, from Lemma 2(3), we get

$$\frac{\|B_2 - B\|_F}{\|B\|_F} = \frac{|x - 1|}{\sqrt{x + 1}} = f_1(x), \quad \frac{\langle B, B_2 \rangle}{\|B\|_F \|B_2\|_F} = \frac{2}{\sqrt{2 + x + 1/x}} = f_2(x),$$

where $x = (0.02c)^s \in (0, 20^{-1/4}]$. It is easy to verify that f_1 decreases monotonically and f_2 increases monotonically for $0 < x < 1$. Hence

$$f_1(x) \geq f_1(20^{-1/4}) \geq 0.4, \quad f_2(x) \leq f_2(20^{-1/4}) \leq 0.94.$$

The proof is completed. \square

Lemma 12. Let $\{X_t\}_{t=1}^{t=T}$ be a sequence of symmetric matrices, and $A_t = \sum_{s=1}^t X_s$, where $t = 1, \dots, T$. Suppose we have two sequences of symmetric matrices $\{Y_t\}_{t=1}^{t=T}$, $\{Z_t\}_{t=0}^{t=T}$, and a sequence real numbers $\{\rho_t\}_{t=0}^{t=T}$ satisfying

$$Y_t = Z_{t-1} + X_t, \quad \rho_t = \rho_{t-1} + \|Y_t - Z_t\|_2, \quad Z_0 = 0, \rho_0 = 0.$$

Define $B_t = \rho_t I + Z_t$, where I denotes the identity matrix. Then for $t = 1, \dots, T$, we have

$$B_t \succeq B_{t-1} + X_t, \quad A_t \preceq B_t \preceq 2\rho_t I + A_t.$$

Proof. Note that for any symmetric matrix \mathbf{S} , it holds that $\|\mathbf{S}\|_2 \mathbf{I} \succeq \mathbf{S}$. Then we have

$$(\rho_t - \rho_{t-1})\mathbf{I} + \mathbf{Z}_t = \|\mathbf{Y}_t - \mathbf{Z}_t\|_2 \mathbf{I} + \mathbf{Z}_t \succeq \mathbf{Y}_t.$$

Adding $\rho_{t-1}\mathbf{I}$ on both sides, we get

$$\mathbf{B}_t = \rho_t \mathbf{I} + \mathbf{Z}_t \succeq \rho_{t-1} \mathbf{I} + \mathbf{Y}_t = \rho_{t-1} \mathbf{I} + \mathbf{Z}_{t-1} + \mathbf{X}_t = \mathbf{B}_{t-1} + \mathbf{X}_t.$$

Hence

$$\mathbf{B}_t = \sum_{s=1}^t (\mathbf{B}_s - \mathbf{B}_{s-1}) \succeq \sum_{s=1}^t \mathbf{X}_s = \mathbf{A}_t.$$

On the other hand, we have

$$\mathbf{Z}_t \preceq \|\mathbf{Z}_t - \mathbf{Y}_t\|_2 \mathbf{I} + \mathbf{Y}_t = (\rho_t - \rho_{t-1})\mathbf{I} + \mathbf{Y}_t.$$

Adding $\rho_t \mathbf{I}$ on both sides, we get

$$\begin{aligned} \mathbf{B}_t &= \rho_t \mathbf{I} + \mathbf{Z}_t \preceq (2\rho_t - \rho_{t-1})\mathbf{I} + \mathbf{Y}_t \\ &= 2(\rho_t - \rho_{t-1})\mathbf{I} + \rho_{t-1}\mathbf{I} + \mathbf{Z}_{t-1} + \mathbf{X}_t \\ &= \mathbf{B}_{t-1} + 2(\rho_t - \rho_{t-1})\mathbf{I} + \mathbf{X}_t. \end{aligned}$$

Hence

$$\mathbf{B}_t = \sum_{s=1}^t (\mathbf{B}_s - \mathbf{B}_{s-1}) \preceq \sum_{s=1}^t 2(\rho_s - \rho_{s-1})\mathbf{I} + \sum_{s=1}^t \mathbf{X}_s = 2\rho_t \mathbf{I} + \mathbf{A}_t.$$

The proof is completed. \square

Theorem 1. Assume that the gradients $\mathbf{G}_1, \dots, \mathbf{G}_T \in \mathbb{R}^{m \times n}$ are matrices of rank at most r . Then for any $\mathbf{W}^* \in \mathbb{R}^{m \times n}$ and $\epsilon > 0$, if $\eta = D/\sqrt{2r}$, the regret of Algorithm 6 is bounded as follows,

$$\sum_{t=1}^T f_t(\mathbf{W}_t) - \sum_{t=1}^T f_t(\mathbf{W}^*) \leq \sqrt{2r}D[2^{1/4}m\rho_T^{1/4} + \text{tr}(\tilde{\mathbf{L}}_T^{1/4})][2^{1/4}n\mu_T^{1/4} + \text{tr}(\tilde{\mathbf{R}}_T^{1/4})],$$

where $D = \max_{t \in [T]} \|\mathbf{W}_t - \mathbf{W}^*\|_F$, $\tilde{\mathbf{L}}_t = \epsilon \mathbf{I}_m + \sum_{s=1}^t \mathbf{G}_s \mathbf{G}_s^\top$, and $\tilde{\mathbf{R}}_t = \epsilon \mathbf{I}_n + \sum_{s=1}^t \mathbf{G}_s^\top \mathbf{G}_s$.

Proof. Define $\hat{\mathbf{L}}_t = (\epsilon + \rho_t)\mathbf{I}_m + \mathbf{L}_t$, $\hat{\mathbf{R}}_t = (\epsilon + \mu_t)\mathbf{I}_n + \mathbf{R}_t$. According to Lemma 12, $\hat{\mathbf{L}}_t$ and $\hat{\mathbf{R}}_t$ are positive definite. Recall the update performed in Algorithm 6,

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta \hat{\mathbf{L}}_t^{-1/4} \mathbf{G}_t \hat{\mathbf{R}}_t^{-1/4}.$$

For $t > 0$, let $\mathbf{H}_t = \hat{\mathbf{L}}_t^{1/4} \otimes \hat{\mathbf{R}}_t^{1/4}$, $\mathbf{g}_t = \overline{\text{vec}}(\mathbf{G}_t)$ and $\mathbf{w}_t = \overline{\text{vec}}(\mathbf{W}_t)$. Due to Lemma 3(3) and Lemma 4, we have

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{H}_t^{-1} \mathbf{g}_t.$$

Lemma 12 implies $0 \prec \hat{\mathbf{L}}_1 \preceq \dots \preceq \hat{\mathbf{L}}_T$, $0 \prec \hat{\mathbf{R}}_1 \preceq \dots \preceq \hat{\mathbf{R}}_T$. Thus, according to Lemma 3(3)(4) and Lemma 6, we get

$$0 \prec \mathbf{H}_1 \preceq \dots \preceq \mathbf{H}_T.$$

Let $\mathbf{H}_0 = \mathbf{0}$. By invoking Lemma 9 and Lemma 8, we obtain the regret bound

$$\begin{aligned} \sum_{t=1}^T f_t(\mathbf{W}_t) - \sum_{t=1}^T f_t(\mathbf{W}^*) &\leq \frac{1}{2\eta} \sum_{t=1}^T (\mathbf{w}_t - \mathbf{w}^*)^\top (\mathbf{H}_t - \mathbf{H}_{t-1}) (\mathbf{w}_t - \mathbf{w}^*) + \frac{\eta}{2} \sum_{t=1}^T (\|\mathbf{g}_t\|_{\mathbf{H}_t}^*)^2 \\ &\leq \frac{D^2}{2\eta} \sum_{t=1}^T \text{tr}(\mathbf{H}_t - \mathbf{H}_{t-1}) + \frac{\eta}{2} \sum_{t=1}^T (\|\mathbf{g}_t\|_{\mathbf{H}_t}^*)^2 \\ &= \frac{D^2}{2\eta} \text{tr}(\mathbf{H}_T) + \frac{\eta}{2} \sum_{t=1}^T (\|\mathbf{g}_t\|_{\mathbf{H}_t}^*)^2, \end{aligned}$$

where $D = \max_{t \in [T]} \|\mathbf{w}_t - \mathbf{w}^*\|_2 = \max_{t \in [T]} \|\mathbf{W}_t - \mathbf{W}^*\|_F$ and $\mathbf{w}^* = \overline{\text{vec}}(\mathbf{W}^*)$.

Define $\widehat{\mathbf{H}}_t = (r\epsilon\mathbf{I} + \sum_{s=1}^t \mathbf{g}_s \mathbf{g}_s^\top)^{1/2}$. Lemma 11 and Lemma 12 imply that

$$\widehat{\mathbf{H}}_t \preceq \sqrt{r} \tilde{\mathbf{L}}_t^{1/4} \otimes \tilde{\mathbf{R}}_t^{1/4} \preceq \sqrt{r} \mathbf{H}_t.$$

Using Lemma 7 and Lemma 10 along with the above equation, we obtain

$$\sum_{t=1}^T (\|\mathbf{g}_t\|_{\mathbf{H}_t}^*)^2 \leq \sqrt{r} \sum_{t=1}^T (\|\mathbf{g}_t\|_{\widehat{\mathbf{H}}_t}^*)^2 \leq 2\sqrt{r} \text{tr}(\widehat{\mathbf{H}}_T) \leq 2r \text{tr}(\mathbf{H}_T).$$

Consequently, using Lemma 3(5) and Lemma 12, we get the desired regret bound

$$\begin{aligned} \sum_{t=1}^T f_t(\mathbf{W}_t) - \sum_{t=1}^T f_t(\mathbf{W}^*) &\leq \left(\frac{D^2}{2\eta} + \eta r \right) \text{tr}(\mathbf{H}_T) = \sqrt{2r} D \text{tr}(\tilde{\mathbf{L}}_T^{1/4}) \text{tr}(\tilde{\mathbf{R}}_T^{1/4}) \\ &\leq \sqrt{2r} D [2^{1/4} m \rho_T^{1/4} + \text{tr}(\tilde{\mathbf{L}}_T^{1/4})] [2^{1/4} n \mu_T^{1/4} + \text{tr}(\tilde{\mathbf{R}}_T^{1/4})], \end{aligned}$$

by choosing $\eta = D/\sqrt{2r}$. The proof is completed. \square

G Experimental Details

We use one RTX3060Ti GPU under the PyTorch 2.0.1+CUDA11.8 framework for DNN training on the CIFAR-100 and Tiny-ImageNet datasets, use one A800 GPU under the PyTorch 2.0.1+CUDA11.7 framework for DNN training on the ImageNet-1k and C4 datasets, and use two NVIDIA L40S GPUs under the PyTorch 2.0.1+CUDA11.8 framework for DNN training on the OWT dataset. To obtain the total peak memory consumption per GPU, we call "torch.cuda.max_memory_allocated".

We set "torch.backends.cudnn.benchmark" to "False" for all the experiments, except when training ViT-Base/32 on the ImageNet-1k dataset. We report the total memory consumption instead of the memory consumption of the second-order optimizer. This total memory includes data, model parameters, activations, gradients, states forming the preconditioners and their inverse roots, states for the used first-order optimizer, and memory fragments. *Our focus lies in quantizing the states for constructing preconditioners and their inverse roots, which are approximately 7x smaller for 4-bit Shampoo compared to 32-bit Shampoo. Because the block size is 64, its maximum value should be calculated every 64 elements and saved as a 32-bit value, resulting in an additional overhead of 0.5 bits (32/64). Consequently, the memory savings are approximately 7 times, calculated as 32/(4+0.5).* In the future, we may adopt double quantization [9] to further reduce memory consumption.

For SGDM, Adagrad or AdamW used in second-order optimizers, we use 32-bit optimizer states on image classification tasks and 16-bit optimizer states on natural language modeling tasks by default. For SGDM, we set the momentum to 0.9 and use an initial learning rate of 0.1. For Adagrad, we set $\epsilon = 10^{-10}$ and use an initial learning rate of 0.01. For AdamW, we set $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\epsilon = 10^{-8}$ and use an initial learning rate of 0.001. For quantization settings, we employ block-wise normalization with a block size of 64 and linear square quantization by default. Matrices with a size smaller than 4096 will not be quantized. For Shampoo and CASPR, we use $\epsilon = 10^{-6}$, $\beta = 0.95$ and $t_1 = 1$, $t_2 = 4$ by default. Shampoo and CASPR precondition blocks from large matrices and the maximum order of a preconditioner is 10000 for 130M LLAMA-2 and is 1200 for other models. For training loss, we use cross-entropy loss. For image classification tasks, automatic mixed precision is enabled except for training transformers on the CIFAR-100 and Tiny-ImageNet datasets.

Settings on training CNNs on CIFAR-100 or Tiny-ImageNet. Minibatch size is set to 128. Weight decay is 0.0005. Data augmentation includes random crop and horizontal flip. For Shampoo, we set $T_1 = 100$ and $T_2 = 500$. In Section 5, we run SGDM for 300 epochs and SGDM+Shampoo for 200 epochs on the CIFAR-100 dataset. We run SGDM for 150 epochs and SGDM+Shampoo for 100 epochs on the Tiny-ImageNet dataset. We adopt the multi-step learning rate schedule (the learning rate is multiplied by 0.1 for every 30% epochs with a linear warmup at the first 5 epochs).

Settings on training transformers on CIFAR-100 or Tiny-ImageNet. We set a patch size of 4 for ViT-small on the CIFAR-100 dataset, and a patch size of 8 for ViT-small on the Tiny-ImageNet dataset. For training Swin-Tiny on the CIFAR-100 dataset, we use a patch size of 2 and window size of 4. For training Swin-Tiny on the Tiny-ImageNet dataset, we use a patch size of 4 and window

size of 7. Minibatch size is set to 128. We run Adagrad/AdamW/NadamW for 150 epochs and Adagrad/AdamW+Shampoo for 100 epochs. Weight decay is 0.0005 for Adagrad, and is 0.05 for AdamW/NadamW. We use the cosine learning rate schedule. Data augmentation follows the source code in [25]. For Shampoo, we set $T_1 = 100$ and $T_2 = 500$. With the exception of certain optimizer settings, the configurations used for ablation studies are identical to those outlined above.

Settings on training ResNet50 on ImageNet-1k. We run SGDM for 120 epochs and SGDM+Shampoo for 100 epochs. Minibatch size is set to 256. Weight decay is 0.0001. We adopt the multi-step learning rate schedule (the learning rate is multiplied by 0.1 for every 30% epochs with a linear warmup at the first 5 epochs). Data augmentation includes random resized crop, horizontal flip, and color jitter. For Shampoo, we set $T_1 = 200$ and $T_2 = 1000$.

Settings on training ViT-Base/32 on ImageNet-1k. We run AdamW for 150 epochs and AdamW+Shampoo for 120 epochs. Minibatch size is set to 512. Weight decay is 0.05. We use the cosine learning rate schedule. Data augmentation follows the configuration for training ViT-Base/16 in [44], excluding repeated augmentation. For Shampoo, we set $T_1 = 200$ and $T_2 = 1000$.

Settings on training GPT-2 on OWT. We run AdamW with 10% warmup steps. Total batch size is set to 480. Batch size is set to 24 for training 124M GPT-2. Dtype is bfloat16. Weight decay is 0.1. For Shampoo, we set $T_1 = 200$ and $T_2 = 200$. For our 4-bit Shampoo, we use Schur-Newton iteration used in Algorithm 4 to compute the inverse root of a preconditioner for training stability.

Settings on training LLAMA-2 on C4. We run AdamW with 10% warmup steps. Total batch size is set to 512. Batch size is set to 256 for training 130M LLAMA-2 and is set to 128 for training 350M LLAMA-2. Dtype is bfloat16. Weight decay is 0. For Shampoo, we set $T_1 = 200$ and $T_2 = 200$.

Settings on K-FAC and AdaBK. K-FAC/AdaBK preconditioners layers without limiting the size of a preconditioner. We set $\beta = 0.9$, $T_1 = 200$, and $T_2 = 2000$. We use $\epsilon = 0.1$ for K-FAC and $\epsilon = 0.001$ for AdaBK. For 4-bit K-FAC/AdaBK, we set $t_1 = 0$ and $t_2 = 0$ (i.e., no orthogonal rectification).

Settings on schedule free optimization. We use the code from [6] to train ResNet34 with SGDScheduleFree and Swin-Tiny with AdamWScheduleFree. For SGDScheduleFree, we set $\text{lr}=1.0$, $\text{weight_decay}=0.0005$ and $\text{warmup_steps}=2000$. For AdamWScheduleFree, we set $\text{lr}=0.0025$, $\text{weight_decay}=0.05$ and $\text{warmup_steps}=10000$.

Settings on M-FAC. We use the code from [15] and set $\text{ngrads}=32$, $\text{damp}=0.1$. The other hyperparameter settings of M-FAC is the same as that of SGDM used for ResNet34 training.

H Additional Results

H.1 Image Classification

More learning rate schedulers. Table 8 shows the performance and wall-clock time of training ResNet34 on CIFAR-100 with cosine learning rate decay. By comparison, SGDM+Shampoo still converges faster than SGDM, and have slightly better test performance.

Table 8: Performance and wall-clock time of training ResNet34 on the CIFAR-100 dataset with cosine learning rate decay. TA = test accuracy, and WCT = wall-clock time.

Epochs	Optimizer	TA (%)	WCT (min)
200	SGDM	79.67	116.0
300	SGDM	79.83	172.7
200	SGDM + 32-bit Shampoo	80.39	152.7
200	SGDM + 4-bit Shampoo (our)	80.22	161.7

We also provide the results of training ResNet34 and Swin-Tiny on CIFAR-100 with schedule-free approach [6] in Table 9. From it one can see that AdamWScheduleFree achieves comparable performance to AdamW with cosine decay, while SGDScheduleFree underperforms compared to SGDM. We observe that this schedule-free algorithm shows rapid improvements in training and test accuracy during the early training stages, but may fail to achieve a higher test accuracy ultimately (see Figure 9). Anyway, these methods are still worse than our AdamW+4-bit Shampoo.

Table 9: Performance and wall-clock time of training on the CIFAR-100 dataset with cosine learning rate decay and schedule-free approach. ResNet34 is trained for 300 epochs and Swin-Tiny is trained for 150 epochs. TA = test accuracy, and WCT = wall-clock time.

Model	Optimizer	TA (%)	WCT (min)
ResNet34	SGDM	79.83	172.7
	SGDScheduleFree	75.63	169.6
Swin-Tiny	AdamW	76.69	318.6
	AdamWScheduleFree	76.58	321.9

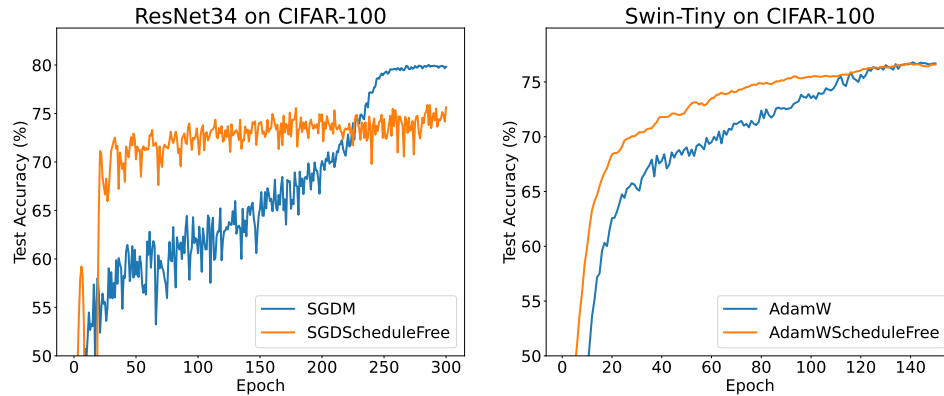


Figure 9: Visualization of test accuracies on the CIFAR-100 dataset with cosine learning rate decay and schedule-free approach.

More optimizers. Table 10 shows results of training Swin-Tiny on CIFAR-100 with NadamW, Adagrad and Adagrad+Shampoo. One can see that Adagrad+4-bit Shampoo converges faster than Adagrad with ignorable extra memory overhead, and also has higher test accuracy. Besides, though NadamW [11] is slightly better than AdamW, it is still worse than our AdamW+4-bit Shampoo.

Table 10: Performance, wall-clock time, and memory cost of training Swin-Tiny on the CIFAR-100 dataset. TA = test accuracy, WCT = wall-clock time, and TMC = total GPU memory cost.

Optimizer	TA (%)	WCT (min)	TMC (MB)
NadamW	77.11	342.4	1465.8
AdamW + 32-bit Shampoo	79.34	260.8	2036.0
AdamW + 4-bit Shampoo (our)	78.63	273.3	1543.9
Adagrad	66.56	294.6	1354.9
Adagrad + 32-bit Shampoo	73.55	245.3	1930.4
Adagrad + 4-bit Shampoo (our)	72.66	259.6	1433.0

M-FAC [15] is a matrix-free method computing inverse-Hessian vector products with many gradient copies. It is not memory-efficient for M-FAC to maintain m dense gradient copies ($m = 1024$ in its official code). Table 11 shows that both SGDM+32-bit Shampoo and SGDM+4-bit Shampoo enjoy much higher efficiency than M-FAC ($m = 32$) for training ResNet34 on CIFAR-100, and enjoy higher test accuracy. EVA [42] is a rank-one second-order optimizer and is memory-efficient. We train ResNet34 on CIFAR-100 with SGDM+EVA, but despite extensive hyper-parameter tuning, we fail to achieve acceleration over SGDM. Instead, we cite EVA’s result of training VGG-19 on CIFAR-100 for 200 epochs (see Table 2 in [42]). The test accuracies of SGDM+EVA and SGDM+Shampoo are 73% and 74.5%, respectively.

Table 11: Performance and memory cost of training ResNet34 on the CIFAR-100 dataset with cosine learning rate decay. All the optimizers are run for 200 epochs. TA = test accuracy, and TMC = total GPU memory cost.

Optimizer	SGDM	M-FAC ($m=32$)	SGDM + 32-bit Shampoo	SGDM + 4-bit Shampoo (our)
TA (%)	79.67	78.56	80.39	80.22
TMC (MB)	822.03	3424.8	1441.8	908.4

Table 12: Performance, wall-clock time, and memory usage per GPU on natural language modeling tasks. VL = validation loss, WCT = wall-clock time, and TMC = total GPU memory cost.

Dataset	Model	Optimizer	VL	WCT (min)	TMC (MB)
C4	LLAMA-130M	AdamW	3.214	346.9	47026
		AdamW + 32-bit Shampoo	3.184	353.7	48813
		AdamW + 4-bit Shampoo (naive)	3.200	353.5	47316
		AdamW + 4-bit Shampoo (our)	3.194	353.1	47318
	LLAMA-350M	AdamW	2.939	2687	54184
		AdamW + 32-bit Shampoo	2.908	2776	59149
		AdamW + 4-bit Shampoo (naive)	2.930	2753	54894
		AdamW + 4-bit Shampoo (our)	2.924	2795	54894
OWT	GPT2-124M	AdamW	2.954	2310	27010
		AdamW + 32-bit Shampoo	2.936	2330	28490
		AdamW + 4-bit Shampoo (naive)	2.953	2359	27209
		AdamW + 4-bit Shampoo (our)	2.944	2311	27209

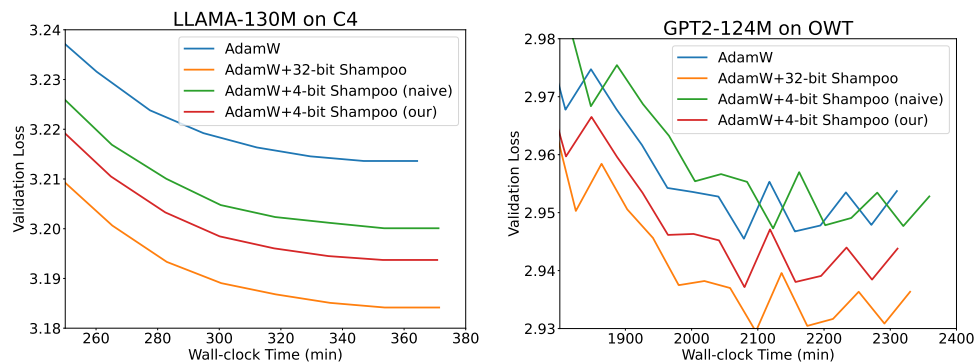


Figure 10: Visualization of validation loss on the C4 and OWT datasets.

H.2 Natural Language Modeling

Models, datasets, and hyperparameters. We train 124M GPT-2 [32] for 60k steps on the OpenWebText (OWT) dataset * following the nanoGPT codebase † with two NVIDIA L40S GPUs, and train 130M LLAMA-2 [37] for 20k steps and 350M LLAMA-2 for 60k steps on the C4 dataset [33] following [43] with one A800 GPU. See Appendix G for experimental details.

Main results. We show the performance, wall-clock time, and memory cost in Table 12, and the validation loss curves in Figure 10. As with the vision tasks, our AdamW+4-bit Shampoo consistently outperformed AdamW and naive AdamW+4-bit Shampoo in terms of performance, and AdamW+32-bit Shampoo in terms of memory usage.

Memory efficiency. We further check the memory usage by increasing token batch size for a language model, which is calculated as the batch size multiplied by the context length (see [43]). To train LLAMA2-7B on the C4 dataset using a single A800 GPU (with a maximum memory of 81,920

*<http://Skylion007.github.io/OpenWebTextCorpus>.

†<https://github.com/karpathy/nanoGPT>.

MB), we set the context length to 256 and then determine the maximum batch size allowed by each optimizer. For Shampoo, the maximum order of a preconditioner for training LLAMA2-7B is 2048. In all experiments, gradient checkpointing is enabled. Table 13 summarizes the evaluation results. By comparison, the 32-bit Shampoo runs out of memory with a batch size of 2, while our 4-bit Shampoo supports a batch size of 64 for standard training and only encounters memory issues at a batch size of 128. These results clearly demonstrate that our 4-bit Shampoo significantly conserves memory compared to the 32-bit version.

Table 13: Memory cost of training LLAMA2-7B on the C4 dataset with different optimizers. One A800 GPU with a maximum memory of 81,920 MB is enabled. TMC = total GPU memory cost, and OOM = out of memory.

Optimizer	Batch Size	TMC (MB)
8-bit AdamW	64	60135
8-bit AdamW	128	68689
8-bit AdamW	256	OOM
8-bit AdamW + 32-bit Shampoo	2	OOM
8-bit AdamW + 4-bit Shampoo (our)	64	74561
8-bit AdamW + 4-bit Shampoo (our)	128	OOM

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: We propose the first second-order optimizers with 4-bit states by taking Shampoo as an example, while preserving the performance achieved with 32-bit optimizer states.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: We discuss the limitations of the work at the end of the paper.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: We provide all the proofs in the Appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.

- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. **Experimental Result Reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We present the implementation details of all the experiments in the Appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. **Open access to data and code**

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We use publicly available datasets and will release our source code.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).

- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. **Experimental Setting/Details**

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [\[Yes\]](#)

Justification: We present the implementation details of all the experiments in the Appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. **Experiment Statistical Significance**

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [\[No\]](#)

Justification: Error bars are not reported because it would be too computationally expensive.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. **Experiments Compute Resources**

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [\[Yes\]](#)

Justification: We present the implementation details of all the experiments in the main paper and Appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.

- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: We review the NeurIPS Code of Ethics and our paper conforms it.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We discuss both potential positive societal impacts and negative societal impacts at the end of the paper.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: We present 4-bit Shampoo for memory efficient training of deep models. It poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.

- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [\[Yes\]](#)

Justification: We properly mention all the existing assets.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [\[Yes\]](#)

Justification: We provide anonymized zip file of our code.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [\[NA\]](#)

Justification: Our paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.