

On the Convergence Analysis of Muon

Wei Shen^{1*}, Ruichuan Huang^{2*}, Minhui Huang³, Cong Shen^{1†}, Jiawei Zhang^{4†}

¹ University of Virginia, ² University of British Columbia
³ Meta, ⁴ University of Wisconsin-Madison

Abstract

The majority of parameters in neural networks are naturally represented as matrices. However, most commonly used optimizers treat these matrix parameters as flattened vectors during optimization, potentially overlooking their inherent structural properties. Recently, an optimizer called Muon has been proposed, specifically designed to optimize matrix-structured parameters. Extensive empirical evidence shows that Muon can significantly outperform traditional optimizers when training neural networks. Nonetheless, the theoretical understanding of Muon’s convergence behavior and the reasons behind its superior performance remain limited. In this work, we present a comprehensive convergence rate analysis of Muon and its comparison with Gradient Descent (GD). We further characterize the conditions under which Muon can outperform GD. Our theoretical results reveal that Muon can benefit from the low-rank and approximate blockwise diagonal structure of Hessian matrices – phenomena widely observed in practical neural network training. Our experimental results support and corroborate the theoretical findings.

1 Introduction

Modern neural networks – such as large language models (LLMs) [Adler et al. \[2024\]](#) – typically consist of a massive amount of parameters and require significant computational resources for training. As a result, designing optimizers that can efficiently train such models has become an important and valuable research problem. Note that most of the parameters in neural networks are naturally represented as matrices – for example, weight matrices in fully connected layers, or the query, key, and value matrices in attention mechanisms. However, most widely adopted optimization algorithms such as Stochastic Gradient Descent (SGD), Adam [\[Diederik, 2014\]](#), and their variants often treat these matrices as flattened vectors, potentially overlooking their inherent structural properties.

Recently, Muon, an optimizer specifically designed for matrix-structured parameters, was proposed in [Jordan et al. \[2024\]](#). The key idea is to update the matrix parameters along an orthogonalized version of their gradients at each iteration (Algorithm 1). Empirical results across various studies and network scales consistently demonstrate the superior performance of Muon in optimizing neural networks with matrix parameters [\[Jordan et al., 2024, Liu et al., 2025a, An et al., 2025, Liu et al., 2025b, Ahn and Xu, 2025\]](#). However, the convergence behavior of Muon and the underlying mechanisms that contribute to its advantage over traditional gradient-descent-based optimizers, such as Gradient Descent (GD) and SGD, are not yet fully understood. Therefore, in this work, we aim to bridge this gap by providing a comprehensive theoretical analysis of Muon’s convergence properties, along with detailed comparisons to GD-based algorithms.

The stepsize for GD usually depends on the Lipschitz smoothness constant (the max singular value of Hessians), which can change rapidly or be very large. This makes it difficult to tune the stepsize. Instead of the max singular value of Hessians (Frobenius norm Lipschitz constant (Assumption 3.1)), by analyzing the convergence behaviors of Muon without assuming the uniform Lipschitz smoothness, we show that the behavior of Muon depends on the average of the global information of the Hessian matrices during training, which can be viewed as smoothing over steps and singular values compared with the Lipschitz constant of GD (the max singular value of Hessians), which gives the potential advantages of Muon over GD.

*Co-first authors

†Co-last authors

Our main contributions can be summarized as follows:

- In the nonconvex setting, we rigorously establish the convergence guarantees for Muon in both deterministic and stochastic cases. We first provide detailed comparisons with GD, demonstrating that Muon can benefit from the low-rank [Sagun et al., 2016, 2017, Wu et al., 2020, Pappas, 2020, Yao et al., 2020] and approximately blockwise diagonal structure of Hessian matrices [Dong et al., 2025, Zhang et al., 2024b,a, Collobert, 2004], which are widely observed phenomena in practical neural network training. Moreover, we analyze the convergence of Muon without assuming the uniform Lipschitz smoothness (Theorem 4.8). We show that the convergence rate of Muon can be related to the average of the global information of the Hessian matrices during training. We characterize the conditions under which Muon outperforms GD, and validate these theoretical findings through experiments on neural network training.
- To further study the advantages of Muon compared to GD, we consider star convex functions (Assumption 4.9), for which we can study the convergence rate of function value. In the star convex setting, under mild assumptions, we show that Muon can achieve a convergence rate comparable to GD’s. Furthermore, we prove that when the Hessian matrices exhibit relatively low-rank and approximately blockwise diagonal structures, Muon can potentially outperform GD. Our experiments on quadratic functions validate and support these theoretical findings.

2 Related Work

Muon was originally proposed in Jordan et al. [2024]. [Jordan et al., 2024] also demonstrated the superior performance of Muon on various models and datasets. Subsequently, Liu et al. [2025a] improved the Muon by incorporating techniques such as weight decay and adjusting the per-parameter update scale, and showed that Muon can outperform Adam when optimizing large language models (LLMs). One possible explanation for Muon’s superior performance is provided by Bernstein and Newhouse [2024], whose work shows that Muon’s update direction corresponds to the steepest descent under the spectral norm constraint. However, during the writing of this paper, we note that there are some concurrent works [Li and Hong, 2025, An et al., 2025, Kovalev, 2025] that also analyze the convergence behavior of Muon. For example, Li and Hong [2025] analyze the Frobenius norm convergence (Definition 3.4) of Muon under Frobenius norm Lipschitz smooth (Assumption 3.1). An et al. [2025] analyze the convergence of the Simplified Muon (Algorithm 2) with Assumption 4.5. Kovalev [2025] provide the convergence analysis of Muon with Assumption 3.2. However, compared to these prior works [Li and Hong, 2025, An et al., 2025, Kovalev, 2025], our work provides a more comprehensive theoretical analysis, i.e. we analyze the convergence of Muon under various smoothness assumptions (Assumption 3.1, 3.2, 4.5), including scenarios without uniform Lipschitz smoothness (Assumption 4.7). Additionally, we offer detailed comparisons with GD and characterize the conditions under which Muon outperforms GD, and validate these conditions through experiments. Therefore, our work presents distinct contributions and insights. We believe that all these works collectively contribute to a deeper understanding of Muon.

Low-rank and approximate blockwise diagonal structure of Hessian. Extensive prior works [Collobert, 2004, Zhang et al., 2024b,a, Dong et al., 2025] have shown that the Hessian of a neural network tends to exhibit a blockwise diagonal structure, with each block corresponding to an individual neuron, both theoretically and empirically. Recently, Zhang et al. [2024b,a] numerically confirmed this property in small Transformers, and Dong et al. [2025] provided theoretical explanations for why such a blockwise diagonal structure emerges in neural network Hessians. In addition, numerous studies [Sagun et al., 2016, 2017, Wu et al., 2020, Pappas, 2020, Yao et al., 2020] have also observed that neural network Hessians are typically low-rank, i.e., their effective ranks, which can be measured by the ratio $\|H\|_*/\|H\|_{\text{op}}$, can be significantly smaller than their dimensionality.

Other structured optimization methods. Although commonly used optimizers for training deep neural networks, such as SGD, Adam [Diederik, 2014], and AdamW [Loshchilov and Hutter, 2017], typically treat structured parameters (e.g., matrices) as flattened vectors, however, in recent years, there has been growing interest in designing structured optimizers that explicitly leverage the inherent structure of parameters. For instance, Adafactor [Shazeer and Stern, 2018], LAMB [You et al., 2019], and Adam-mini [Zhang et al., 2024b] incorporate matrix- or layer-level structure to reduce memory consumption. KFAC [Martens and Grosse, 2015] and TNT [Ren and Goldfarb, 2021]

approximate the Fisher matrix to implement natural gradient methods. Shampoo [Gupta et al., 2018] is specifically designed for matrix or tensor parameters and can be viewed as an approximation to the full-matrix preconditioner in AdaGrad [Duchi et al., 2011]. Morwani et al. [2024] provided additional theoretical analyses and modifications to Shampoo. Recently, SOAP [Vyas et al., 2024] was introduced, which combines the ideas of Adam and Shampoo. Galore [Zhao et al., 2024] was proposed to exploit the low-rank structure of gradients for memory efficiency. Additionally, Liu et al. [2025b] proposed COSMOS, which can be seen as a combination of SOAP and Muon. More recently, Xie et al. [2025] and An et al. [2025] introduced ASGO (One-Sided Shampoo), which only applies a one-sided preconditioner within the Shampoo. Xie et al. [2025], An et al. [2025] showed that ASGO(One-Sided Shampoo) can achieve better convergence rate than the original Shampoo.

3 Preliminaries

Notation. For a vector v , we denote its l_2 norm as $\|v\|_2$. For a matrix $A \in \mathbb{R}^{m \times n}$, we denote its nuclear norm as $\|A\|_*$, spectral norm as $\|A\|_{\text{op}}$, Frobenius norm as $\|A\|_F$, and Λ norm as $\|A\|_\Lambda = \sqrt{\text{tr}(A\Lambda A^\top)}$ where $\Lambda \in \mathbb{R}^{n \times n}$ is a real symmetric positive definite matrix. For $A \in \mathbb{R}^{m \times n}$ with rows a_1, \dots, a_m , we define $\text{vec}(A) = (a_1^\top a_2^\top \dots a_m^\top)^\top \in \mathbb{R}^{mn}$ as the vectorization of A . We denote I_d as the identity matrix with d dimension. For matrix $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{m' \times n'}$, we denote their Kronecker product as $A \otimes B \in \mathbb{R}^{mm' \times nn'}$. We define the operator norm of a third-order tensor $H \in \mathbb{R}^{d_1 \times d_2 \times d_3}$ as $\|H\|_{\text{op}} = \max_{\|u\|_2=\|v\|_2=\|w\|_2=1} |\sum_{i=1}^{d_1} \sum_{j=1}^{d_2} \sum_{k=1}^{d_3} H_{ijk} u_i v_j w_k|$, where $u \in \mathbb{R}^{d_1}, v \in \mathbb{R}^{d_2}, w \in \mathbb{R}^{d_3}$.

In this paper, we consider the following stochastic optimization problem:

$$\min_{W \in \mathbb{R}^{m \times n}} f(W) = \mathbb{E}_\xi[f(W; \xi)]. \quad (1)$$

We denote $f^* = \inf_W f(W)$, $r = \min\{m, n\}$, $f_v(w) = f(W)$, and $w = \text{vec}(W) \in \mathbb{R}^{mn}$. We assume $f^* > -\infty$. We consider the following smoothness assumptions in this paper.

Assumption 3.1 (Frobenius norm Lipschitz smooth). We say $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ is L Frobenius norm Lipschitz smooth if for any $W, W' \in \mathbb{R}^{m \times n}$, we have

$$\|\nabla f(W) - \nabla f(W')\|_F \leq L \|W - W'\|_F.$$

Assumption 3.2 (Spectral norm Lipschitz smooth). We say $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ is L_* spectral norm Lipschitz smooth if for any $W, W' \in \mathbb{R}^{m \times n}$, we have

$$\|\nabla f(W) - \nabla f(W')\|_* \leq L_* \|W - W'\|_{\text{op}}.$$

The first is the standard Frobenius norm-based smoothness Assumption 3.1, which serves as a natural extension of the conventional l_2 -norm smoothness for functions with vector parameters to functions with matrix parameters. The second is a spectral norm smoothness Assumption 3.2, which accounts for the distinct structure of matrix parameters. In this work, we establish the convergence of Muon under both assumptions.

For a stochastic setting, we also use the following standard bounded variance assumption.

Assumption 3.3 (Bounded variance). We assume $\nabla f(W; \xi)$ is an unbiased stochastic estimator of the true gradient $\nabla f(W)$ and have a bounded variance, i.e.

$$\mathbb{E}[\nabla f(W; \xi)] = \nabla f(W), \quad \mathbb{E}\|\nabla f(W; \xi) - \nabla f(W)\|_F^2 \leq \sigma^2.$$

Definition 3.4. We say \widehat{W} is an ϵ -Frobenius norm stationary point of f if $\|\nabla f(\widehat{W})\|_F \leq \epsilon$.

Definition 3.5. We say \widehat{W} is an ϵ -nuclear norm stationary point of f if $\|\nabla f(\widehat{W})\|_* \leq \epsilon$.

For a matrix A with rank r , there is a well-known relationship between its Frobenius norm and nuclear norm: $\|A\|_F \leq \|A\|_* \leq \sqrt{r} \|A\|_F$. Thus, we have following proposition.

Proposition 3.6. If $\widehat{W} \in \mathbb{R}^{m \times n}$ is an ϵ -nuclear norm stationary point of f , then it is also an ϵ -Frobenius norm stationary point of f . If $\widehat{W} \in \mathbb{R}^{m \times n}$ is an ϵ -Frobenius norm stationary point of f , then it is also an $\sqrt{r}\epsilon$ -nuclear norm stationary point of f , where $r = \min\{m, n\}$.

We note that for functions with matrix parameters, the nuclear norm stationarity is tighter than the ordinary Frobenius norm stationarity.

3.1 Muon

Algorithm 1 Muon

```

1: Input: Initial weights  $W_0$ , learning rate schedule  $\{\eta_t\}$ ,  $\beta \in [0, 1)$ , batch size  $B$ 
2: for  $t = 0, 1, \dots, T - 1$  do
3:   Sample batch  $\{\xi_{t,i}\}_{i=1}^B$  uniformly
4:    $G_t = \frac{1}{B} \sum_{i=1}^B \nabla f(W_t; \xi_{t,i})$  (or  $G_t = \nabla f(W_t)$  in the deterministic setting)
5:   If  $t > 0$ ,  $M_t = \beta M_{t-1} + (1 - \beta)G_t$ . If  $t = 0$ ,  $M_0 = G_0$ .
6:    $(U_t, S_t, V_t) = \text{SVD}(M_t)$ 
7:    $W_{t+1} = W_t - \eta_t U_t V_t^\top$ 
8: end for

```

Algorithm 2 Simplified Muon

```

1: Input: Initial weights  $W_0$ , learning rate schedule  $\{\eta_t\}$ 
2: for  $t = 0, 1, \dots, T - 1$  do
3:    $G_t = \nabla f(W_t)$ 
4:    $(U_t, S_t, V_t) = \text{SVD}(G_t)$ 
5:    $W_{t+1} = W_t - \eta_t U_t V_t^\top$ 
6: end for

```

In this subsection, we introduce the Muon algorithm. The original idea of Muon in [Jordan et al. \[2024\]](#) is to orthogonalize the update matrix. For example, suppose the original update direction is G_t with rank r_t and the singular value decomposition (SVD) of G_t is $U_t S_t V_t^\top$, where $S_t \in \mathbb{R}^{r_t \times r_t}$ is the diagonal matrix of the singular values of G_t , $U_t \in \mathbb{R}^{m \times r_t}$ and $V_t \in \mathbb{R}^{n \times r_t}$ are the left and right singular vector matrices of G_t , respectively. Then, in Muon, the update matrix will be $U_t V_t^\top$, which is the nearest semi-orthogonal matrix to the original G_t . Following this main idea, we have [Algorithm 2](#), which can be viewed as the simplest form of Muon.

Additionally, if we add momentum and conduct orthogonalization over the momentum direction, we have [Algorithm 1](#), which is the main algorithm we will analyze in the stochastic setting. In fact, our later experimental and theoretical results show that [Algorithm 1](#) and [2](#) have very similar performance in the deterministic setting, while additional momentum can be helpful in the stochastic setting (See [Figure 1](#)).

Furthermore, in practice, computing the SVD is very costly. Therefore, a common practical implementation of Muon uses Newton–Schulz iterations [Bernstein and Newhouse \[2024\]](#), [Jordan et al. \[2024\]](#) to approximate the orthogonalization process. An introduction of the Muon with Newton–Schulz iterations can be found in [Appendix E](#). However, our later experiments show that the performance of the SVD-based version of Muon is similar to the Newton–Schulz-based version, with the primary difference being the higher computational cost of the SVD procedure. (See [Figure 1](#)). Therefore, without loss of generality, we primarily analyze the convergence properties of [Algorithm 1](#) and [2](#) using SVD in this paper.

4 Convergence of Muon

4.1 Nonconvex case

In this subsection, we analyze the convergence guarantees of Muon in the nonconvex setting with various smoothness assumptions. All the proofs of theorems in this subsection can be found in [Appendix B](#). We first prove the convergence of Muon with the Frobenius norm Lipschitz smoothness assumption.

Theorem 4.1. *Under Assumptions [3.1](#) and [3.3](#), if we apply [Algorithm 1](#) with $\eta_t = \eta$, then*

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(W_t)\|_*] \leq \frac{\mathbb{E}[f(W_0) - f(W_T)]}{T\eta} + \frac{Lr\eta}{2} + \frac{2\sigma\sqrt{r(1-\beta)}}{\sqrt{B(1+\beta)}} + \frac{2\beta\sigma\sqrt{r}}{(1-\beta)T\sqrt{B}} + \frac{2r\eta\beta L}{1-\beta}.$$

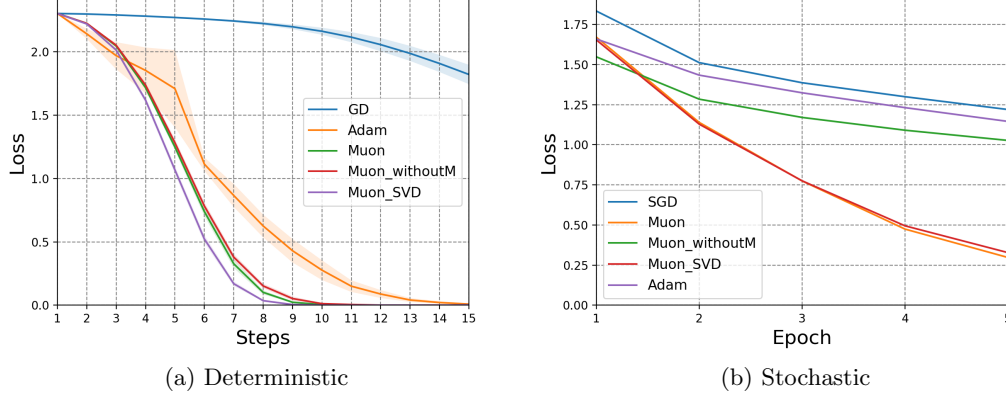


Figure 1: Comparison of (S)GD, Adam, Muon (using Newton–Schulz iterations with momentum), Muon_withoutM (using Newton–Schulz iterations without momentum), and Muon_SVD (Algorithm 1). In the deterministic setting (a), loss is defined and trained over a fixed subset of CIFAR-10 [Krizhevsky et al., 2009]. In the stochastic setting (b), training utilizes mini-batches randomly sampled from the complete CIFAR-10 training set. The loss is evaluated on the entire CIFAR-10 training set.

Denote $\Delta = f(W_0) - f^*$. If we set $B = 1$, $\eta = \sqrt{\frac{(1-\beta)\Delta}{rTL}}$, $1 - \beta = \min\{\frac{\sqrt{L\Delta}}{\sigma\sqrt{T}}, 1\}$, then

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(W_t)\|_*] \leq O\left(\sqrt[4]{\frac{r^2 L \Delta \sigma^2}{T}} + \sqrt{\frac{r L \Delta}{T}} + \frac{\sqrt{r} \sigma^2}{\sqrt{L \Delta T}}\right).$$

Thus, we can find an ϵ -nuclear norm stationary point of f with a complexity of $O(r^2 L \sigma^2 \Delta \epsilon^{-4})$.

Note that if we set $\beta = O(1)$ (or $\beta = 0$) in the stochastic setting, a large batch size $B = O(\epsilon^{-2})$ is required to sufficiently reduce the variance of the stochastic mini-batch gradient. Therefore, our theoretical analysis suggests that momentum can be beneficial in the stochastic setting, which aligns with the observation in Figure 1 where Muon with momentum outperforms Muon without momentum in stochastic setting.

We have the following corollary for the deterministic setting.

Corollary 4.2. Under Assumption 3.1, if we apply Algorithm 1 with $G_t = \nabla f(W_t)$ and $\eta_t = \eta$, then

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(W_t)\|_* \leq \frac{f(W_0) - f(W_T)}{T\eta} + \frac{Lr\eta}{2} + \frac{2r\eta\beta L}{1-\beta}$$

Denote $\Delta = f(W_0) - f^*$. If we set $\eta = \sqrt{\frac{\Delta}{rTL}}$, $\beta = O(1)$, then

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(W_t)\|_* \leq O\left(\sqrt{\frac{rL\Delta}{T}}\right).$$

Thus, we can find an ϵ -nuclear norm stationary point of f with a complexity of $O(rL\Delta\epsilon^{-2})$.

Our theoretical results indicate that in the deterministic setting, setting $\beta = O(1)$ (or even $\beta = 0$) does not significantly affect the convergence guarantee. This is consistent with the empirical observation in Figure 1, where Muon with momentum performs similarly to Muon without momentum in the deterministic setting.

Note that the (S)GD's complexity of finding ϵ -Frobenius norm stationary points of f is $O(L\sigma^2\Delta\epsilon^{-4})$ in the stochastic setting and $O(L\Delta\epsilon^{-2})$ in the deterministic setting [Garrigos and Gower, 2023], which means (S)GD can also find ϵ -nuclear norm stationary points of f with $O(r^2 L \sigma^2 \Delta \epsilon^{-4})$ complexity in the stochastic setting and $O(rL\Delta\epsilon^{-2})$ in the deterministic setting according to Proposition 3.6. Thus, there is little direct advantage of Muon over (S)GD when analyzing the convergence rate with the conventional Frobenius norm Lipschitz smoothness assumption.

However, if we adopt the spectral norm Lipschitz smoothness assumption, the comparison changes.

Theorem 4.3. Under Assumptions 3.2 and 3.3, if we apply Algorithm 1 with $\eta_t = \eta$, then

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(W_t)\|_*] \leq \frac{\mathbb{E}[f(W_0) - f(W_T)]}{T\eta} + \frac{L_*\eta}{2} + \frac{2\sigma\sqrt{r(1-\beta)}}{\sqrt{(1+\beta)B}} + \frac{2\beta\sigma\sqrt{r}}{(1-\beta)T\sqrt{B}} + \frac{2\eta\beta L_*}{1-\beta}.$$

Denote $\Delta = f(W_0) - f^*$. If we set $B = 1$, $\eta = \sqrt{\frac{(1-\beta)\Delta}{TL_*}}$, $1 - \beta = \min\{\frac{\sqrt{L_*\Delta}}{\sigma\sqrt{rT}}, 1\}$, then

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(W_t)\|_*] \leq O\left(\sqrt[4]{\frac{rL_*\Delta\sigma^2}{T}} + \sqrt{\frac{L_*\Delta}{T}} + \frac{r\sigma^2}{\sqrt{L_*\Delta T}}\right).$$

Thus, we can find an ϵ -nuclear norm stationary point of f with a complexity of $O(rL_*\sigma^2\Delta\epsilon^{-4})$.

Similarly, we have the following corollary for the deterministic setting.

Corollary 4.4. Under Assumption 3.2, if we apply Algorithm 1 with $G_t = \nabla f(W_t)$ and $\eta_t = \eta$, then

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(W_t)\|_* \leq \frac{f(W_0) - f(W_T)}{T\eta} + \frac{L_*\eta}{2} + \frac{2\eta\beta L_*}{1-\beta}.$$

Denote $\Delta = f(W_0) - f^*$. If we set $\eta = \sqrt{\frac{\Delta}{TL_*}}$, $\beta = O(1)$, then

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(W_t)\|_* \leq O\left(\sqrt{\frac{L_*\Delta}{T}}\right).$$

Thus, we can find an ϵ -nuclear norm stationary point of f with a complexity of $O(L_*\Delta\epsilon^{-2})$.

For a function f , its L and L_* are related to its Hessians $\nabla^2 f$. In neural network optimization, a widely observed phenomenon is that Hessians are typically low-rank [Sagun et al., 2016, 2017, Wu et al., 2020, Pappas, 2020, Yao et al., 2020] and are approximately blockwise diagonal [Dong et al., 2025, Zhang et al., 2024b,a, Collobert, 2004] (e.g., See Figure 2 in Zhang et al. [2024a] and Figure 1-3 in Dong et al. [2025]). Thus, we make the following assumption to better illustrate the connections and differences between L and L_* in neural networks.

Assumption 4.5 (Λ -norm Lipschitz smooth). We say $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ is 1 - Λ -norm Lipschitz smooth if for any $W, W' \in \mathbb{R}^{m \times n}$,

$$\|\nabla f(W) - \nabla f(W')\|_{\Lambda^{-1}} \leq \|W - W'\|_{\Lambda}.$$

where $\Lambda \in \mathbb{R}^{n \times n}$ is a positive definite matrix.

Assumption 4.5 is closely related to the blockwise diagonal structure of the Hessian and is equivalent to the the following assumption of Hessians:

$$-I_m \otimes \Lambda = - \begin{bmatrix} \Lambda & & \\ & \Lambda & \\ & & \dots \\ & & & \Lambda \end{bmatrix} \preceq \nabla^2 f_v(w) \preceq \begin{bmatrix} \Lambda & & \\ & \Lambda & \\ & & \dots \\ & & & \Lambda \end{bmatrix} = I_m \otimes \Lambda \in \mathbb{R}^{mn \times mn}, \quad (2)$$

where $f_v(w) = f(W)$ and $w = \text{vec}(W) \in \mathbb{R}^{mn}$. Assumption 4.5 has also been used in An et al. [2025] for analyzing matrix-structured optimization.

Lemma 4.6. If $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ is 1 - Λ -norm Lipschitz smooth with a positive definite matrix $\Lambda \in \mathbb{R}^{n \times n}$, then f is also $\|\Lambda\|_{\text{op}}$ Frobenius norm Lipschitz smooth and $\|\Lambda\|_*$ spectral norm Lipschitz smooth.

Thus, if the Hessians of f have a blockwise diagonal structure and are relatively “low rank”, i.e. if it satisfies Assumption 4.5 with Λ such that $\|\Lambda\|_{\text{op}} \approx \|\Lambda\|_* \ll r\|\Lambda\|_{\text{op}}$, then the convergence rate of Muon for finding the nuclear norm stationary points are better than (S)GD’s, i.e. $O(\|\Lambda\|_*\Delta\epsilon^{-2}) \ll O(r\|\Lambda\|_{\text{op}}\Delta\epsilon^{-2})$ and $O(r\|\Lambda\|_*\sigma^2\Delta\epsilon^{-4}) \ll O(r^2\|\Lambda\|_{\text{op}}\sigma^2\Delta\epsilon^{-4})$ for deterministic and stochastic settings respectively.

To further investigate how the structure and dynamics of the Hessian matrix affect the convergence of Muon and GD in a more fine-grained manner, we consider the following assumption and analysis.

Assumption 4.7. We assume the norm of $\nabla^3 f$ is bounded, i.e., we assume for any $W \in \mathbb{R}^{m \times n}$, $\|\nabla^3 f_v(w)\|_{\text{op}} \leq s$. Here $f_v(w) = f(W)$ and $w = \text{vec}(W) \in \mathbb{R}^{mn}$.

As discussed before, the momentum term is mainly used to reduce the variance of the stochastic mini-batch gradient. It plays a similar role as the standard SGD with momentum. The key difference between Muon and GD is the modified update direction – the orthogonal update direction. Thus, for simplicity, we consider Algorithm 2, the Muon without momentum in the deterministic setting.

Theorem 4.8. Under Assumption 4.7, if we apply Algorithm 2 with $\eta_t = \eta$, then

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(W_t)\|_* \leq \frac{f(W_0) - f(W_T)}{T\eta} + \frac{\eta J}{2} + \frac{s\eta^2 r^{3/2}}{6}$$

where $J = \frac{1}{T} \sum_{t=0}^{T-1} J_t$, $J_t = \text{vec}(U_t V_t^\top)^\top H_t \text{vec}(U_t V_t^\top)$, $H_t = \nabla^2 f_v(\text{vec}(W_t))$, and $f_v(\text{vec}(W_t)) = f(W_t)$. Denote $\Delta = f(W_0) - f^*$. If $J > 0$ ¹, and $\eta = \sqrt{\frac{2\Delta}{JT}}$, then

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(W_t)\|_* \leq \sqrt{\frac{2J\Delta}{T}} + \frac{sr^{3/2}\Delta}{3JT}.$$

Thus, we can find an ϵ -nuclear norm stationary point of f with a complexity of $O(J\Delta\epsilon^{-2})$.

Note that

$$\begin{aligned} J_t &= \text{vec}(U_t V_t^\top)^\top H_t \text{vec}(U_t V_t^\top) \\ &= \text{vec}(I_{r_t})^\top (U_t^\top \otimes V_t^\top) H_t (U_t \otimes V_t) \text{vec}(I_{r_t}) \\ &= \sum_{i,j \in [r_t]} [(U_t^\top \otimes V_t^\top) H_t (U_t \otimes V_t)]_{(i-1)*r_t+i, (j-1)*r_t+j}. \end{aligned} \quad (3)$$

Thus, J_t is the sum of r_t^2 elements of $A_t \triangleq (U_t^\top \otimes V_t^\top) H_t (U_t \otimes V_t)$, where $A_t \in \mathbb{R}^{r_t^2 \times r_t^2}$ can be viewed as a representation of H_t under a certain congruence transformation depending on U_t and V_t . In fact, if we assume $H_t = P_t \otimes Q_t$, where $P_t \in \mathbb{R}^{m \times m}$ and $Q_t \in \mathbb{R}^{n \times n}$, then $A_t = (U_t^\top P_t U_t) \otimes (V_t^\top Q_t V_t)$ and $J_t = \langle U_t^\top P_t U_t, V_t^\top Q_t V_t \rangle \leq \|U_t^\top P_t U_t\|_{\text{op}} \|V_t^\top Q_t V_t\|_* \leq \|P_t\|_{\text{op}} \|Q_t\|_*$. Note that the complexity of GD for ϵ -nuclear norm stationary point is $O(rL\Delta\epsilon^{-2})$, where L is at least larger than $\max_t L_t$ with $L_t \triangleq \|H_t\|_{\text{op}} = \max_t \|P_t\|_{\text{op}} \|Q_t\|_{\text{op}}$. Thus, when H_t can be represented by $P_t \otimes Q_t$ and Q_t (or P_t) are relatively low-rank such that $\|Q_t\|_{\text{op}} \approx \|Q_t\|_* \ll r \|Q_t\|_{\text{op}}$, $J = \frac{1}{T} \sum_t J_t \leq \frac{1}{T} \sum_t \|P_t\|_{\text{op}} \|Q_t\|_* \ll r \max_t L_t \leq rL$, and the convergence rate of Muon can be better than GD's.

However, various studies [Gur-Ari et al., 2018, Zhao et al., 2021, 2024, Cosson et al., 2023] have shown that the gradients during neural network optimization are also typically low rank. Thus, when we convert the stationary point of Frobenius norm to nuclear norm, the additional coefficient can be much smaller than \sqrt{r} (as shown in Figure 2 (c,g)). Therefore, to compare the convergence rates of Muon and GD more precisely and more fairly, one should examine the ratios between the nuclear norm and Frobenius norm of their gradients as well as J_t and L_t . For example, if, for $t = 0, 1, \dots, T$,

$$\frac{J_t}{L_t} \leq \frac{\|\nabla f(W_t)\|_*^2}{\|\nabla f(W_t)\|_F^2}, \quad (4)$$

we can claim Muon is better than GD for finding the stationary point of f .

We consider investigating and validating (4) in the optimization process of neural network training. For example, we consider the optimization process of a three-layer Multi-Layer Perceptron (MLP) for an image classification task on the MNIST dataset [LeCun et al., 1998]. We optimize this neural network using GD and Muon (Algorithm 2) respectively, and record the loss, J_t , L_t , $\|\nabla f(W_t)\|_*^2$ and $\|\nabla f(W_t)\|_F^2$ (with respect to one of the matrix parameters) during the optimization process. The results are shown in Figure 2. Detailed settings can be found in Appendix D. From Figure 2, we can observe that in the early stages of optimization, J_t can actually be smaller than L_t (b, f), while

¹We find that usually $J > 0$ in experiments. Thus, we mainly discuss the situations when $J > 0$ in the main paper. When $J < 0$, there is a better convergence rate regarding T ; see discussions in Appendix B.

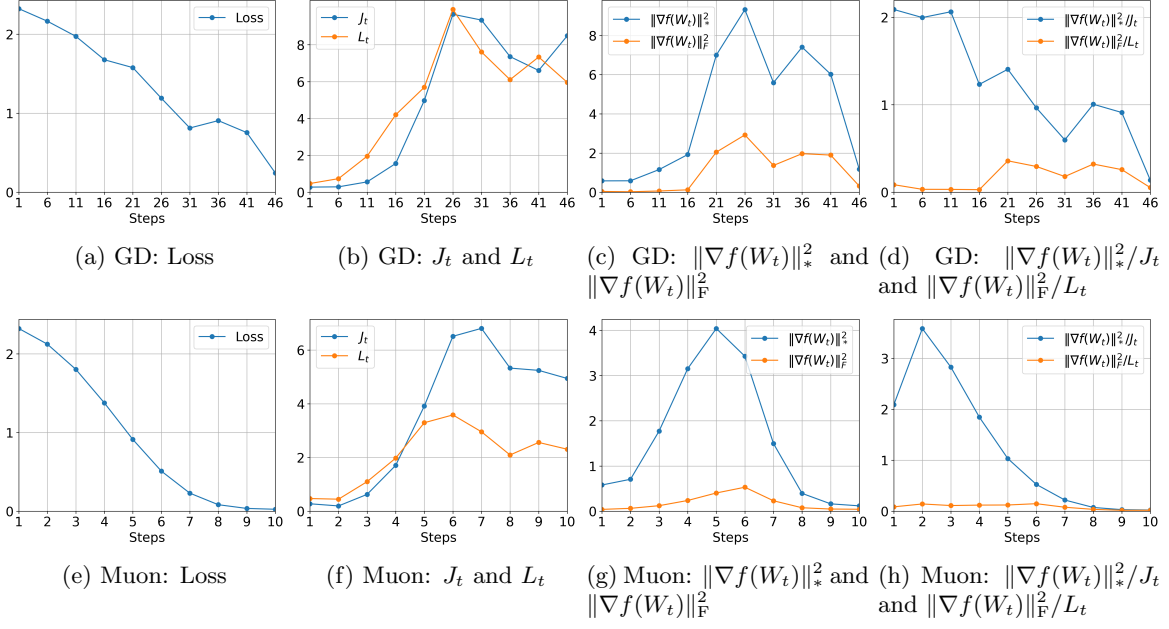


Figure 2: Comparison of J_t , L_t , $\|\nabla f(W_t)\|_*^2$ and $\|\nabla f(W_t)\|_F^2$ over the training process of GD and Muon (Algorithm 2). f is defined as the cross-entropy loss of a MLP with three matrix parameters $W^1 \in \mathbb{R}^{128 \times 784}$, $W^2 \in \mathbb{R}^{64 \times 128}$, $W^3 \in \mathbb{R}^{10 \times 64}$ over a fixed subset of MNIST. We show the gradients and Hessians with respect to W^2 in this Figure. Detailed settings can be found in Appendix D.

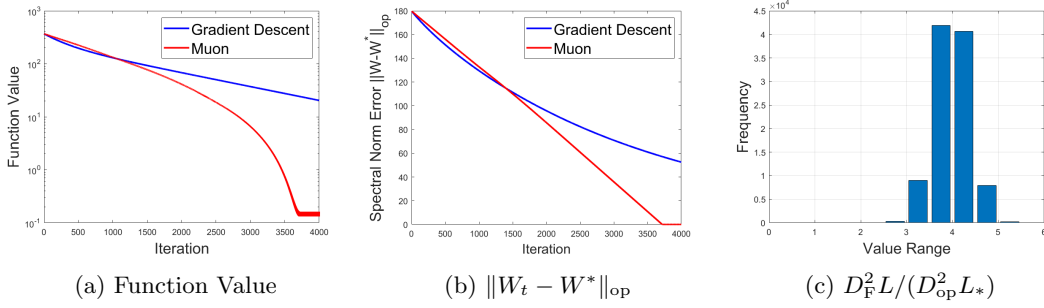


Figure 3: Experiments on a quadratic function. Detailed settings can be found in Appendix D.

$\|\nabla f(W_t)\|_*^2$ is significantly larger than $\|\nabla f(W_t)\|_F^2$ (d, g). As a result, Muon’s ratio ($\|\nabla f(W_t)\|_*^2 / J_t$) can be significantly larger than GD’s ($\|\nabla f(W_t)\|_F^2 / L_t$) (d, h). Therefore, according to our theory in Equation (4), Muon can converge much faster than GD in the early phase, which is also corroborated in our experimental results (a, e). As the optimization approaches convergence, J_t becomes comparable to or even larger than L_t (b, f). However, throughout the entire optimization process of both GD and Muon, the ratio of Muon ($\|\nabla f(W_t)\|_*^2 / J_t$) remains consistently larger than that of GD ($\|\nabla f(W_t)\|_F^2 / L_t$) (d, h), demonstrating the advantage of Muon over GD in optimizing neural network matrix parameters.

4.2 Star-Convex case

To further study the advantages of Muon, we consider the star convex case so that we can study the convergence of function value instead of the gradient norm.

Assumption 4.9 (Star convex). For any $W \in \mathbb{R}^{m \times n}$, we assume the following inequality holds

$$\langle \nabla f(W), W - W^* \rangle \geq f(W) - f^*.$$

where $W^* \in \operatorname{argmin}_W f(W)$ is the optimal point.

Note that the star convex condition [Nesterov and Polyak, 2006] is weaker than the convex condition, and empirical evidence [Kleinberg et al., 2018, Zhou et al., 2019] suggests that, in some cases, the optimization path of neural networks can satisfy this condition.

Moreover, to ensure the convergence of Muon, we adopt the following mild assumption. Note that the same assumption is also used in Gupta et al. [2018], An et al. [2025] for their convergence analysis of their algorithms in similar settings.

Assumption 4.10. For $W_t, t = 0, \dots, T$, generated by Algorithm 2, we assume $\|W_t - W^*\|_{\text{op}} \leq D_{\text{op}}$.

From Figure 3(b), we can observe that when we apply Muon on a quadratic function, $\|W_t - W^*\|$ is bounded, which aligns with Assumption 4.10.

With the spectral norm Lipschitz smoothness assumption, we have the following convergence theorem of Muon.

Theorem 4.11. Under Assumption 4.9, 3.2 and 4.10, if we apply Algorithm 2 with $\eta_t = \eta$, then

$$f(W_T) - f^* \leq \left(1 - \frac{\eta}{D_{\text{op}}}\right)^T (f(W_0) - f^*) + \frac{L_* D_{\text{op}} \eta}{2}.$$

Thus, to reach the precision $f(W_T) - f^* \leq \epsilon$, we can take $\eta = O(\frac{\epsilon}{L_* D_{\text{op}}})$, and the complexity is $O(L_* D_{\text{op}}^2 \epsilon^{-1} \log \frac{\Delta}{\epsilon})$, where $\Delta = f(W_0) - f^*$.

Moreover, if we apply Algorithm 2 with adaptive learning rates $\eta_t = \frac{\|\nabla f(W_t)\|_*}{L_*}$, we have

$$f(W_t) - f^* \leq \frac{2L_*(f(W_0) - f^*)D_{\text{op}}^2}{2L_* D_{\text{op}}^2 + t(f(W_0) - f^*)}.$$

Thus, to reach the precision $f(W_T) - f^* \leq \epsilon$, the complexity is $O(L_* D_{\text{op}}^2 \epsilon^{-1})$.

The proof of Theorem 4.11 can be found in Appendix C. Moreover, similar to the last section, we also conduct convergence analysis of Muon under Assumption 3.1 and 4.7 in Appendix C.

Recall that, in the analysis of GD, to reach the precision $f(W_T) - f^* \leq \epsilon$, the complexity is $O(LD_{\text{F}}^2 \epsilon^{-1})$, where L is the smoothness constant under Frobenius norm and $D_{\text{F}} = \|W_0 - W^*\|_{\text{F}}$ [Garrigos and Gower, 2023]. Hence, to compare the convergence rate of Muon and GD, we need to compare LD_{F}^2 with $L_* D_{\text{op}}^2$. Similar to the discussions in the last section, when the Hessians of f have a blockwise diagonal structure and are relatively “low rank”, i.e. if it satisfies Assumption 4.5 with Λ such that $\|\Lambda\|_{\text{op}} \approx \|\Lambda\|_* \ll r\|\Lambda\|_{\text{op}}$, then $L_* D_{\text{op}}^2 \leq \|\Lambda\|_* D_{\text{op}}^2 \ll r\|\Lambda\|_{\text{op}} D_{\text{op}}^2$. In experiments, we usually find that $\|W_t - W^*\|_{\text{op}}$ are decreasing as the t increase (Figure 3(b)). Thus, we assume $D_{\text{op}} \approx \|W_0 - W^*\|_{\text{op}}$. Though, the success of LoRA [Hu et al., 2022] in fine-tuning foundation models indicating a relatively low rank of $W_0 - W^*$ is enough for fine-tuning, however, many evidences [Lialin et al., Jiang et al., 2024, Huang et al., 2025] demonstrate that a relatively high rank of $W_0 - W^*$ is needed for pretraining and some complex fine-tuning tasks for large foundation models, low rank of $W_0 - W^*$ can lead to degraded performance in those scenarios. Thus, we assume $W_0 - W^*$ is relatively “high rank”. Then, we can expect $r\|W_0 - W^*\|_{\text{op}}^2 \approx \|W_0 - W^*\|_{\text{F}}^2$ and $L_* D_{\text{op}}^2 \ll r\|\Lambda\|_{\text{op}} D_{\text{op}}^2 \approx LD_{\text{F}}^2$, which means Muon can have a better performance than GD.

In Figure 3, we conduct experiments comparing Muon and GD on quadratic function $f(W) = \text{tr}((W - W^*)^\top Q(W - W^*))$, where Q is a positive definite matrix with a poor condition number and is relatively “low rank”, i.e. $\|Q\|_* \approx \|Q\|_{\text{op}}$. Note that the Hessian of f is $I \otimes Q$, which satisfies Assumption 4.5 and is aligned with our pervious discussions. We randomly choose W^* , and set $W_0 = 0$. We can note that Muon outperforms GD in Figure 3 and most samples of ratio $D_{\text{F}}^2 L / (D_{\text{op}}^2 L_*)$ are larger than 2, which validates the theoretical findings.

5 Conclusions and future directions

In this work, we established the convergence guarantees of Muon in both nonconvex and star convex settings. We provided detailed analysis and comparisons between Muon and GD, and established the conditions under which Muon can outperform GD. We showed that Muon can exploit the low-rank and approximately blockwise diagonal structure of Hessian matrices – a property commonly observed

during neural network training. Our experiments on neural networks and quadratic functions supported and corroborated the theoretical findings. One of the limitations of our work is that our experiments are conducted on relatively small-scale neural networks, and one direction for future research is to investigate the structure of Hessian matrices in larger-scale neural networks. Moreover, studying the structures of Hessian matrices from theoretical perspectives and leveraging their structures to design new or better optimization methods are also interesting future directions.

References

- Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, et al. Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>, 2:6, 2024.
- Kwangjun Ahn and Byron Xu. Dion: A communication-efficient optimizer for large models. *arXiv preprint arXiv:2504.05295*, 2025.
- Kang An, Yuxing Liu, Rui Pan, Shiqian Ma, Donald Goldfarb, and Tong Zhang. Asgo: Adaptive structured gradient optimization. *arXiv preprint arXiv:2503.20762*, 2025.
- Jeremy Bernstein and Laker Newhouse. Old optimizer, new norm: An anthology. *arXiv preprint arXiv:2409.20325*, 2024.
- Ronan Collobert. Large scale machine learning. 2004.
- Romain Cosson, Ali Jadbabaie, Anuran Makur, Amirhossein Reisizadeh, and Devavrat Shah. Low-rank gradient descent. *IEEE Open Journal of Control Systems*, 2:380–395, 2023.
- Kingma Diederik. Adam: A method for stochastic optimization. (*No Title*), 2014.
- Zhaorui Dong, Yushun Zhang, Zhi-Quan Luo, Jianfeng Yao, and Ruoyu Sun. Towards quantifying the hessian structure of neural networks. *arXiv preprint arXiv:2505.02809*, 2025.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- Guillaume Garrigos and Robert M Gower. Handbook of convergence theorems for (stochastic) gradient methods. *arXiv preprint arXiv:2301.11235*, 2023.
- Vineet Gupta, Tomer Koren, and Yoram Singer. Shampoo: Preconditioned stochastic tensor optimization. In *International Conference on Machine Learning*, pages 1842–1850. PMLR, 2018.
- Guy Gur-Ari, Daniel A Roberts, and Ethan Dyer. Gradient descent happens in a tiny subspace. *arXiv preprint arXiv:1812.04754*, 2018.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- Qiushi Huang, Tom Ko, Zhan Zhuang, Lilian Tang, and Yu Zhang. Hira: Parameter-efficient hadamard high-rank adaptation for large language models. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Ting Jiang, Shaohan Huang, Shengyue Luo, Zihan Zhang, Haizhen Huang, Furu Wei, Weiwei Deng, Feng Sun, Qi Zhang, Deqing Wang, et al. Mora: High-rank updating for parameter-efficient fine-tuning. *arXiv preprint arXiv:2405.12130*, 2024.
- Keller Jordan, Yuchen Jin, Vlado Boza, Jiacheng You, Franz Cesista, Laker Newhouse, and Jeremy Bernstein. Muon: An optimizer for hidden layers in neural networks, 2024. URL <https://kellerjordan.github.io/posts/muon/>.
- Bobby Kleinberg, Yuanzhi Li, and Yang Yuan. An alternative view: When does sgd escape local minima? In *International conference on machine learning*, pages 2698–2707. PMLR, 2018.

- Dmitry Kovalev. Understanding gradient orthogonalization for deep learning via non-euclidean trust-region optimization. *arXiv preprint arXiv:2503.12645*, 2025.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.(2009), 2009.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Jiaxiang Li and Mingyi Hong. A note on the convergence of muon and further. *arXiv preprint arXiv:2502.02900*, 2025.
- Vladislav Lialin, Namrata Shivagunde, Sherin Muckatira, and Anna Rumshisky. Relora: High-rank training through low-rank updates, 2023. URL <https://arxiv.org/abs/2307.05695>.
- Jingyuan Liu, Jianlin Su, Xingcheng Yao, Zhejun Jiang, Guokun Lai, Yulun Du, Yidao Qin, Weixin Xu, Enzhe Lu, Junjie Yan, et al. Muon is scalable for llm training. *arXiv preprint arXiv:2502.16982*, 2025a.
- Liming Liu, Zhenghao Xu, Zixuan Zhang, Hao Kang, Zichong Li, Chen Liang, Weizhu Chen, and Tuo Zhao. Cosmos: A hybrid adaptive optimizer for memory-efficient training of llms. *arXiv preprint arXiv:2502.17410*, 2025b.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pages 2408–2417. PMLR, 2015.
- Depen Morwani, Itai Shapira, Nikhil Vyas, Eran Malach, Sham Kakade, and Lucas Janson. A new perspective on shampoo’s preconditioner. *arXiv preprint arXiv:2406.17748*, 2024.
- Yurii Nesterov and Boris T Polyak. Cubic regularization of newton method and its global performance. *Mathematical programming*, 108(1):177–205, 2006.
- Vardan Papyan. Traces of class/cross-class structure pervade deep learning spectra. *Journal of Machine Learning Research*, 21(252):1–64, 2020.
- Yi Ren and Donald Goldfarb. Tensor normal training for deep learning models. *Advances in Neural Information Processing Systems*, 34:26040–26052, 2021.
- Levent Sagun, Leon Bottou, and Yann LeCun. Eigenvalues of the hessian in deep learning: Singularity and beyond. *arXiv preprint arXiv:1611.07476*, 2016.
- Levent Sagun, Utku Evci, V Ugur Guney, Yann Dauphin, and Leon Bottou. Empirical analysis of the hessian of over-parametrized neural networks. *arXiv preprint arXiv:1706.04454*, 2017.
- Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, pages 4596–4604. PMLR, 2018.
- Nikhil Vyas, Depen Morwani, Rosie Zhao, Mujin Kwun, Itai Shapira, David Brandfonbrener, Lucas Janson, and Sham Kakade. Soap: Improving and stabilizing shampoo using adam. *arXiv preprint arXiv:2409.11321*, 2024.
- Yikai Wu, Xingyu Zhu, Chenwei Wu, Annie Wang, and Rong Ge. Dissecting hessian: Understanding common structure of hessian in neural networks. *arXiv preprint arXiv:2010.04261*, 2020.
- Shuo Xie, Tianhao Wang, Sashank Reddi, Sanjiv Kumar, and Zhiyuan Li. Structured preconditioners in adaptive optimization: A unified analysis. *arXiv preprint arXiv:2503.10537*, 2025.
- Zhewei Yao, Amir Gholami, Kurt Keutzer, and Michael W Mahoney. Pyhessian: Neural networks through the lens of the hessian. In *2020 IEEE international conference on big data (Big data)*, pages 581–590. IEEE, 2020.

- Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep learning: Training bert in 76 minutes. *arXiv preprint arXiv:1904.00962*, 2019.
- Yushun Zhang, Congliang Chen, Tian Ding, Ziniu Li, Ruoyu Sun, and Zhiquan Luo. Why transformers need adam: A hessian perspective. *Advances in Neural Information Processing Systems*, 37:131786–131823, 2024a.
- Yushun Zhang, Congliang Chen, Ziniu Li, Tian Ding, Chenwei Wu, Diederik P Kingma, Yinyu Ye, Zhi-Quan Luo, and Ruoyu Sun. Adam-mini: Use fewer learning rates to gain more. *arXiv preprint arXiv:2406.16793*, 2024b.
- Jiawei Zhao, Florian Tobias Schaefer, and Anima Anandkumar. Zero initialization: Initializing residual networks with only zeros and ones. 2021.
- Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong Tian. Galore: Memory-efficient llm training by gradient low-rank projection. *arXiv preprint arXiv:2403.03507*, 2024.
- Yi Zhou, Junjie Yang, Huishuai Zhang, Yingbin Liang, and Vahid Tarokh. Sgd converges to global minimum in deep learning via star-convex path. *arXiv preprint arXiv:1901.00451*, 2019.

Appendix

The Appendix is organized as follows. In Appendix A, we introduce some lemmas that will be utilized in the subsequent proofs and present the proof of Lemma 4.6. In Appendix B, we present the proofs of theorems in the nonconvex setting. In Appendix C, we present the proofs of theorems in the star convex setting. In Appendix D, we provide the detailed settings of our experiments. Moreover, in Appendix E, we provide an introduction to the Newton–Schulz iteration, which is commonly used in the practical implementation of Muon, and show how it approximates the orthogonalization process.

A Lemmas

Lemma A.1 (Lemma 8 in An et al. [2025]). *For a symmetric positive definite matrix $\Lambda \in \mathbb{R}^{n \times n}$ and matrix $A \in \mathbb{R}^{m \times n}$, it holds that*

$$\|A\|_* \leq \sqrt{\|\Lambda\|_*} \|A\|_{\Lambda^{-1}}.$$

Lemma A.2. *For a symmetric positive definite matrix $\Lambda \in \mathbb{R}^{n \times n}$ and matrix $A \in \mathbb{R}^{m \times n}$, it holds that*

$$\|A\|_F \leq \sqrt{\|\Lambda\|_{\text{op}}} \|A\|_{\Lambda^{-1}}, \quad (5)$$

$$\|A\|_{\Lambda} \leq \sqrt{\|\Lambda\|_{\text{op}}} \|A\|_F, \quad (6)$$

$$\|A\|_{\Lambda} \leq \sqrt{\|\Lambda\|_*} \|A\|_{\text{op}}. \quad (7)$$

Proof. We first prove Equation (5) and Equation (6).

Suppose the spectral decomposition of Λ is $\Lambda = U \Sigma U^\top$, where U is an orthogonal matrix, $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_n)$, $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$ are the singular values of Λ . Denote $AU = B = (b_1 \ b_2 \ \dots \ b_n)$, where b_i is the i -th column of B . Then, we have

$$\|A\|_{\Lambda}^2 = \text{tr}(A \Lambda A^\top) = \text{tr}(A U \Sigma U^\top A^\top) = \text{tr}(\Sigma B^\top B) = \sum_{i=1}^n \sigma_i \|b_i\|_2^2.$$

Note that

$$\|A\|_F^2 = \sum_{i=1}^n \|b_i\|_2^2.$$

Thus, we have

$$\|A\|_{\Lambda} \leq \sqrt{\sigma_1} \|A\|_F = \sqrt{\|\Lambda\|_{\text{op}}} \|A\|_F,$$

and

$$\|A\|_{\Lambda} \geq \sqrt{\sigma_n} \|A\|_F.$$

Note that $\sigma_n^{-1} = \|\Lambda^{-1}\|_{\text{op}}$. Thus,

$$\|A\|_F \leq \sqrt{\sigma_n^{-1}} \|A\|_{\Lambda} = \|\Lambda^{-1}\|_{\text{op}} \|A\|_{\Lambda}. \quad (8)$$

Since Equation (8) holds for any Λ , we can define $\Lambda' = \Lambda^{-1}$ and get

$$\|A\|_F \leq \|\Lambda'\|_{\text{op}} \|A\|_{\Lambda'^{-1}}. \quad (9)$$

Thus, we proved Equation (5) and Equation (6). For Equation (7), we can prove it with the following inequalities

$$\|A\|_{\Lambda} = \sqrt{\text{tr}(A \Lambda A^\top)} = \sqrt{\text{tr}(A^\top A \Lambda)} \leq \sqrt{\|A^\top A\|_{\text{op}} \|\Lambda\|_*} = \sqrt{\|\Lambda\|_*} \|A\|_{\text{op}}.$$

□

Proof of Lemma 4.6

Proof. If $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ is 1 - Λ -norm Lipschitz smooth with a positive definite matrix $\Lambda \in \mathbb{R}^{n \times n}$, then, for any $W, W' \in \mathbb{R}^{m \times n}$, we have

$$\begin{aligned}\|\nabla f(W) - \nabla f(W')\|_F &\leq \sqrt{\|\Lambda\|_{\text{op}}} \|\nabla f(W) - \nabla f(W')\|_{\Lambda^{-1}} \\ &\leq \sqrt{\|\Lambda\|_{\text{op}}} \|W - W'\|_{\Lambda} \\ &\leq \|\Lambda\|_{\text{op}} \|W - W'\|_F\end{aligned}$$

where the first inequality is due to Equation (5) in Lemma A.2, the second inequality is the definition of the Λ -smoothness, and the third inequality is due to Equation (6) in Lemma A.2. Thus, f is also $\|\Lambda\|_{\text{op}}$ Frobenius norm Lipschitz smooth.

Moreover, we have

$$\begin{aligned}\|\nabla f(W) - \nabla f(W')\|_* &\leq \sqrt{\|\Lambda\|_*} \|\nabla f(W) - \nabla f(W')\|_{\Lambda^{-1}} \\ &\leq \sqrt{\|\Lambda\|_*} \|W - W'\|_{\Lambda} \\ &\leq \|\Lambda\|_* \|W - W'\|_{\text{op}}\end{aligned}$$

where the first inequality is due to Lemma A.1, the second inequality is the definition of the Λ -smoothness, and the third inequality is due to Equation (7) in Lemma A.2. Thus, f is also $\|\Lambda\|_*$ spectral norm Lipschitz smooth. \square

Lemma A.3. For $t = 0, 1, \dots, T$, M_t and W_t generated by Algorithm 1, defining $C_0 = \nabla f(W_0)$, and when $t > 0$, $C_t = \beta C_{t-1} + (1 - \beta) \nabla f(W_t) = (1 - \beta) \sum_{i=1}^t \beta^{t-i} \nabla f(W_i) + \beta^t \nabla f(W_0)$, we have

$$\mathbb{E}[\|C_t - M_t\|_F] \leq \sqrt{\frac{1 - \beta}{1 + \beta}} \frac{\sigma}{\sqrt{B}} + \beta^t \frac{\sigma}{\sqrt{B}}.$$

Proof. According to Assumption 3.3, we have following relationship about G_t and $\nabla f(W_t)$.

$$\begin{aligned}\mathbb{E}[\|G_t - \nabla f(W_t)\|_F^2] &= \mathbb{E} \left[\left\| \frac{1}{B} \sum_{i=1}^B \nabla f(W_t; \xi_{t,i}) - \nabla f(W_t) \right\|_F^2 \right] \\ &= \frac{1}{B^2} \sum_{i=1}^B \mathbb{E}[\|\nabla f(W_t; \xi_{t,i}) - \nabla f(W_t)\|_F^2] \\ &\leq \frac{\sigma^2}{B}\end{aligned}\tag{10}$$

where the second equality is due to $\mathbb{E}[\nabla f(W_t; \xi_{t,i})] = \nabla f(W_t)$ and the last inequality is due to $\mathbb{E}[\|\nabla f(W; \xi_{t,i}) - \nabla f(W)\|_F^2] \leq \sigma^2$. Moreover, according to the Cauchy-Schwarz inequality, we have

$$\mathbb{E}[\|G_t - \nabla f(W_t)\|_F] \leq \sqrt{\mathbb{E}[\|G_t - \nabla f(W_t)\|_F^2]} = \frac{\sigma}{\sqrt{B}}.\tag{11}$$

Note that $M_0 = G_0$, $M_t = \beta M_{t-1} + (1 - \beta) G_t = (1 - \beta) \sum_{i=1}^t \beta^{t-i} G_i + \beta^t G_0$. Thus, we have

$$\begin{aligned}\mathbb{E}[\|C_t - M_t\|_F] &\leq (1 - \beta) \mathbb{E}[\| \sum_{i=1}^t \beta^{t-i} (G_i - \nabla f(W_i)) \|_F] + \beta^t \mathbb{E}[\|G_0 - \nabla f(W_0)\|_F] \\ &\leq \sqrt{(1 - \beta)^2 \mathbb{E}[\| \sum_{i=1}^t \beta^{t-i} (G_i - \nabla f(W_i)) \|_F^2]} + \beta^t \frac{\sigma}{\sqrt{B}} \\ &\leq \sqrt{(1 - \beta)^2 \sum_{i=1}^t \beta^{t-i} \frac{\sigma^2}{B}} + \beta^t \frac{\sigma}{\sqrt{B}}\end{aligned}$$

$$\leq \sqrt{\frac{1-\beta}{1+\beta}} \frac{\sigma}{\sqrt{B}} + \beta^t \frac{\sigma}{\sqrt{B}}$$

where the second inequality is due to the Cauchy-Schwarz inequality and Equation (11), the third inequality is due to Equation (10). \square

B Nonconvex

B.1 Proof of Theorem 4.1

Proof. Set $\eta_t = \eta$. Since f is L Frobenius norm Lipschitz smooth, we have

$$\begin{aligned} \mathbb{E}[f(W_t) - f(W_{t+1})] &\geq \mathbb{E}[\eta \langle \nabla f(W_t), U_t V_t^\top \rangle - \frac{L}{2} \eta^2 \|U_t V_t^\top\|_F^2] \\ &\geq \mathbb{E}[\eta \langle M_t, U_t V_t^\top \rangle - \frac{rL}{2} \eta^2 - \eta \langle \nabla f(W_t) - M_t, U_t V_t^\top \rangle] \\ &\geq \mathbb{E}[\eta \langle M_t, U_t V_t^\top \rangle - \frac{rL}{2} \eta^2 - \eta \|\nabla f(W_t) - M_t\|_F \|U_t V_t^\top\|_F] \\ &\geq \mathbb{E}[\eta \|M_t\|_* - \frac{rL}{2} \eta^2 - \eta \sqrt{r} \|\nabla f(W_t) - M_t\|_F] \\ &\geq \mathbb{E}[\eta \|\nabla f(W_t)\|_* - \frac{rL}{2} \eta^2 - \eta(\sqrt{r} + 1) \|\nabla f(W_t) - M_t\|_F] \\ &\geq \mathbb{E}[\eta \|\nabla f(W_t)\|_* - \frac{rL}{2} \eta^2 - 2\eta \sqrt{r} \|\nabla f(W_t) - M_t\|_F]. \end{aligned}$$

Then, we need to analyze and bound the error $\|\nabla f(W_t) - M_t\|_F$. Note that $M_0 = G_0$, $M_t = \beta M_{t-1} + (1 - \beta)G_t = (1 - \beta) \sum_{i=1}^t \beta^{t-i} G_i + \beta^t G_0$. We can define $C_0 = \nabla f(W_0)$, and when $t > 0$, $C_t = \beta C_{t-1} + (1 - \beta) \nabla f(W_t) = (1 - \beta) \sum_{i=1}^t \beta^{t-i} \nabla f(W_i) + \beta^t \nabla f(W_0)$. Then, according to Lemma A.3, we have

$$\mathbb{E}[\|C_t - M_t\|_F] \leq \sqrt{\frac{1-\beta}{1+\beta}} \frac{\sigma}{\sqrt{B}} + \beta^t \frac{\sigma}{\sqrt{B}}.$$

Moreover, when $t > 0$, we can note that

$$\begin{aligned} &\mathbb{E}[\|\nabla f(W_t) - C_t\|_F] \\ &= \mathbb{E}[\|\nabla f(W_t) - (\beta C_{t-1} + (1 - \beta) \nabla f(W_t))\|_F] \\ &= \mathbb{E}[\beta \|\nabla f(W_t) - C_{t-1}\|_F] \\ &\leq \mathbb{E}[\beta \|\nabla f(W_{t-1}) - C_{t-1}\|_F + \beta \|\nabla f(W_{t-1}) - \nabla f(W_t)\|_F] \\ &\leq \mathbb{E}[\beta \|\nabla f(W_{t-1}) - C_{t-1}\|_F + \beta L \|W_{t-1} - W_t\|_F] \\ &= \mathbb{E}[\beta \|\nabla f(W_{t-1}) - C_{t-1}\|_F + \beta L \eta \|U_{t-1} V_{t-1}^\top\|_F] \\ &\leq \mathbb{E}[\beta \|\nabla f(W_{t-1}) - C_{t-1}\|_F + \beta L \eta \sqrt{r}] \\ &\leq \beta^t \|\nabla f(W_0) - C_0\|_F + \sum_{i=1}^t \beta^i L \eta \sqrt{r} \\ &\leq \frac{\sqrt{r} \beta L \eta}{1 - \beta} \end{aligned}$$

where the second inequality is due to Assumption 3.1.

Thus, we have

$$\mathbb{E}[\|\nabla f(W_t) - M_t\|_F] \leq \mathbb{E}[\|C_t - M_t\|_F + \|\nabla f(W_t) - C_t\|_F] \leq \sqrt{\frac{1-\beta}{1+\beta}} \frac{\sigma}{\sqrt{B}} + \beta^t \frac{\sigma}{\sqrt{B}} + \frac{\sqrt{r} \beta L \eta}{1 - \beta},$$

and

$$\mathbb{E}[f(W_t) - f(W_{t+1})] \geq \mathbb{E}[\eta \|\nabla f(W_t)\|_* - \frac{rL}{2} \eta^2 - 2\eta \sqrt{r} \|\nabla f(W_t) - M_t\|_F]$$

$$\geq \mathbb{E}[\eta \|\nabla f(W_t)\|_*] - \frac{rL}{2}\eta^2 - 2\eta \sqrt{\frac{1-\beta}{1+\beta}} \frac{\sigma\sqrt{r}}{\sqrt{B}} - 2\eta\beta^t \frac{\sigma\sqrt{r}}{\sqrt{B}} - \frac{2r\eta^2\beta L}{1-\beta}.$$

Summing over $t = 0, 1, \dots, T-1$, we can get

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(W_t)\|_*] \leq \frac{1}{T\eta} \mathbb{E}[f(W_0) - f(W_T)] + \frac{Lr\eta}{2} + \frac{2\sigma\sqrt{r(1-\beta)}}{\sqrt{(1+\beta)B}} + \frac{2\beta\sigma\sqrt{r}}{(1-\beta)T\sqrt{B}} + \frac{2r\eta\beta L}{1-\beta}.$$

Denote $\Delta = f(W_0) - f^*$ and set $\eta = \sqrt{\frac{(1-\beta)\Delta}{rTL}}$, we can get

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(W_t)\|_*] \leq \frac{\Delta}{T\eta} + \frac{Lr\eta}{2} + \frac{2\sigma\sqrt{r(1-\beta)}}{\sqrt{(1+\beta)B}} + \frac{2\beta\sigma\sqrt{r}}{(1-\beta)T\sqrt{B}} + \frac{2r\eta\beta L}{1-\beta}.$$

When $B = 1$, we can set $1 - \beta = \min\{\frac{\sqrt{L\Delta}}{\sigma\sqrt{T}}, 1\}$ and get

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(W_t)\|_*] \leq O(1) \left(\sqrt[4]{\frac{r^2 L \Delta \sigma^2}{T}} + \sqrt{\frac{r L \Delta}{T}} + \frac{\sqrt{r} \sigma^2}{\sqrt{L \Delta T}} \right).$$

Thus, we can find an ϵ -nuclear norm stationary point of f with a complexity of $O(r^2 L \sigma^2 \Delta \epsilon^{-4})$. \square

B.2 Proof of Theorem 4.3

Proof. Set $\eta_t = \eta$. Since f is L_* spectral norm Lipschitz smooth, we have

$$\begin{aligned} \mathbb{E}[f(W_t) - f(W_{t+1})] &\geq \mathbb{E}[\eta \langle \nabla f(W_t), U_t V_t^\top \rangle - \frac{L_*}{2} \eta^2 \|U_t V_t^\top\|_{\text{op}}^2] \\ &\geq \mathbb{E}[\eta \langle M_t, U_t V_t^\top \rangle - \frac{L_*}{2} \eta^2 - \eta \langle \nabla f(W_t) - M_t, U_t V_t^\top \rangle] \\ &\geq \mathbb{E}[\eta \langle M_t, U_t V_t^\top \rangle - \frac{L_*}{2} \eta^2 - \eta \|\nabla f(W_t) - M_t\|_* \|U_t V_t^\top\|_{\text{op}}] \\ &\geq \mathbb{E}[\eta \|M_t\|_* - \frac{L_*}{2} \eta^2 - \eta \|\nabla f(W_t) - M_t\|_*] \\ &\geq \mathbb{E}[\eta \|\nabla f(W_t)\|_* - \frac{L_*}{2} \eta^2 - 2\eta \|\nabla f(W_t) - M_t\|_*]. \end{aligned}$$

Then, we need to analyze and bound the error $\|\nabla f(W_t) - M_t\|_*$. Note that $M_0 = G_0$, $M_t = \beta M_{t-1} + (1-\beta)G_t = (1-\beta) \sum_{i=1}^t \beta^{t-i} G_i + \beta^t G_0$. We can define $C_0 = \nabla f(W_0)$, and when $t > 0$, $C_t = \beta C_{t-1} + (1-\beta) \nabla f(W_t) = (1-\beta) \sum_{i=1}^t \beta^{t-i} \nabla f(W_i) + \beta^t \nabla f(W_0)$. Then, according to Lemma A.3, we have

$$\mathbb{E}[\|C_t - M_t\|_*] \leq \sqrt{r} \mathbb{E}[\|C_t - M_t\|_{\text{F}}] \leq \sqrt{\frac{1-\beta}{1+\beta}} \frac{\sigma\sqrt{r}}{\sqrt{B}} + \beta^t \frac{\sigma\sqrt{r}}{\sqrt{B}}.$$

Moreover, when $t > 0$, we can note that

$$\begin{aligned} &\mathbb{E}[\|\nabla f(W_t) - C_t\|_*] \\ &= \mathbb{E}[\|\nabla f(W_t) - (\beta C_{t-1} + (1-\beta) \nabla f(W_t))\|_*] \\ &= \mathbb{E}[\beta \|\nabla f(W_t) - C_{t-1}\|_*] \\ &\leq \mathbb{E}[\beta \|\nabla f(W_{t-1}) - C_{t-1}\|_* + \beta \|\nabla f(W_{t-1}) - \nabla f(W_t)\|_*] \\ &\leq \mathbb{E}[\beta \|\nabla f(W_{t-1}) - C_{t-1}\|_* + \beta L_* \|W_{t-1} - W_t\|_{\text{op}}] \\ &= \mathbb{E}[\beta \|\nabla f(W_{t-1}) - C_{t-1}\|_* + \beta L_* \eta \|U_{t-1} V_{t-1}^\top\|_{\text{op}}] \\ &\leq \mathbb{E}[\beta \|\nabla f(W_{t-1}) - C_{t-1}\|_* + \beta L_* \eta] \\ &\leq \beta^t \|\nabla f(W_0) - C_0\|_{\text{F}} + \sum_{i=1}^t \beta^i L_* \eta \end{aligned}$$

$$\leq \frac{\beta L_* \eta}{1 - \beta}$$

where the second inequality is due to Assumption 3.2.

Thus, we have

$$\mathbb{E}[\|\nabla f(W_t) - M_t\|_*] \leq \mathbb{E}[\|C_t - M_t\|_* + \|\nabla f(W_t) - C_t\|_*] \leq \sqrt{\frac{1 - \beta}{1 + \beta}} \frac{\sigma \sqrt{r}}{\sqrt{B}} + \beta^t \frac{\sigma \sqrt{r}}{\sqrt{B}} + \frac{\beta L_* \eta}{1 - \beta},$$

and

$$\begin{aligned} \mathbb{E}[f(W_t) - f(W_{t+1})] &\geq \mathbb{E}[\eta \|\nabla f(W_t)\|_* - \frac{L_*}{2} \eta^2 - 2\eta \|\nabla f(W_t) - M_t\|_*] \\ &\geq \mathbb{E}[\eta \|\nabla f(W_t)\|_* - \frac{L_*}{2} \eta^2 - 2\eta \sqrt{\frac{1 - \beta}{1 + \beta}} \frac{\sigma \sqrt{r}}{\sqrt{B}} - 2\eta \beta^t \frac{\sigma \sqrt{r}}{\sqrt{B}} - \frac{2\eta^2 \beta L_*}{1 - \beta}]. \end{aligned}$$

Summing over $t = 0, 1, \dots, T - 1$, we can get

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(W_t)\|_*] \leq \frac{1}{T\eta} \mathbb{E}[f(W_0) - f(W_T)] + \frac{L_* \eta}{2} + \frac{2\sigma \sqrt{r(1 - \beta)}}{\sqrt{(1 + \beta)B}} + \frac{2\beta \sigma \sqrt{r}}{(1 - \beta)T\sqrt{B}} + \frac{2\eta \beta L_*}{1 - \beta}.$$

Denote $\Delta = f(W_0) - f^*$ and set $\eta = \sqrt{\frac{(1 - \beta)\Delta}{TL_*}}$, we can get

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(W_t)\|_*] \leq \frac{\Delta}{T\eta} + \frac{L_* \eta}{2} + \frac{2\sigma \sqrt{r(1 - \beta)}}{\sqrt{(1 + \beta)B}} + \frac{2\beta \sigma \sqrt{r}}{(1 - \beta)T\sqrt{B}} + \frac{2\eta \beta L_*}{1 - \beta}.$$

When $B = 1$, we can set $1 - \beta = \min\{\frac{\sqrt{L_* \Delta}}{\sigma \sqrt{rT}}, 1\}$ and get

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(W_t)\|_*] \leq O(1) \left(\sqrt[4]{\frac{r L_* \Delta \sigma^2}{T}} + \sqrt{\frac{L_* \Delta}{T}} + \frac{r \sigma^2}{\sqrt{L_* \Delta T}} \right).$$

Thus, we can find an ϵ -nuclear norm stationary point of f with a complexity of $O(r L_* \sigma^2 \Delta \epsilon^{-4})$. \square

B.3 Proof of Theorem 4.8

Proof. Denote $d = mn$, $J_t = \text{vec}(U_t V_t^\top)^\top H_t \text{vec}(U_t V_t^\top)$, $H_t = \nabla^2 f_v(\text{vec}(W_t))$, and $f_v(\text{vec}(W_t)) = f(W_t)$. Set $\eta_t = \eta$. Taking the Taylor expansion at W_t , we have

$$\begin{aligned} f(W_{t+1}) - f(W_t) &= \langle \nabla f(W_t), W_{t+1} - W_t \rangle + \frac{1}{2} \text{vec}(W_{t+1} - W_t)^\top \nabla^2 f_v(\text{vec}(W_t)) \text{vec}(W_{t+1} - W_t) \\ &\quad + \frac{1}{6} \sum_{i,j,k=1}^d [\nabla^3 f_v(\theta)]_{ijk} \text{vec}(W_{t+1} - W_t)_i \text{vec}(W_{t+1} - W_t)_j \text{vec}(W_{t+1} - W_t)_k \\ &\leq -\eta \langle \nabla f(W_t), U_t V_t^\top \rangle + \frac{\eta^2 J_t}{2} + \frac{s \eta^3 \|\text{vec}(U_t V_t^\top)\|_F^3}{6} \\ &\leq -\eta \langle \nabla f(W_t), U_t V_t^\top \rangle + \frac{\eta^2 J_t}{2} + \frac{s \eta^3 r^{3/2}}{6} \\ &= -\eta \|\nabla f(W_t)\|_* + \frac{\eta^2 J_t}{2} + \frac{s \eta^3 r^{3/2}}{6} \end{aligned} \tag{12}$$

where $\theta \in \mathbb{R}^d$ can be some vectors between $\text{vec}(W_t)$ and $\text{vec}(W_{t+1})$, and the first inequality is due to Assumption 4.7.

Denote $J = \frac{1}{T} \sum_{t=0}^{T-1} J_t$. Summing over $t = 0, 1, \dots, T - 1$, we can get

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(W_t)\|_* \leq \frac{f(W_0) - f(W_T)}{T\eta} + \frac{\eta J}{2} + \frac{s \eta^2 r^{3/2}}{6}.$$

Denote $\Delta = f(W_0) - f^*$. In experiments, we found that usually $J > 0$. When $J > 0$ and set $\eta = \sqrt{\frac{2\Delta}{JT}}$, we can get

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(W_t)\|_* \leq \sqrt{\frac{2J\Delta}{T}} + \frac{sr^{3/2}\Delta}{3JT}.$$

Thus, when $J > 0$, we can find an ϵ -nuclear norm stationary point of f with a complexity of $O(J\Delta\epsilon^{-2})$.

When $J \leq 0$, we have

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(W_t)\|_* \leq \frac{\Delta}{T\eta} + \frac{s\eta^2 r^{3/2}}{6}.$$

We can set $\eta = \sqrt[3]{\frac{\Delta}{sTr^{3/2}}}$ and get

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(W_t)\|_* \leq O\left(\frac{\Delta^{2/3} s^{1/3} r^{1/2}}{T^{2/3}}\right).$$

Thus, when $J \leq 0$, we can find an ϵ -nuclear norm stationary point of f with a complexity of $O(s^{1/2} r^{3/4} \Delta \epsilon^{-3/2})$. \square

C Star-Convex

In this section, we first present the proof of Theorem 4.11. Then, we present an additional convergence analysis of Muon under Assumption 3.1 and Assumption 4.7 in Appendix C.3 and Appendix C.2 respectively.

C.1 Proof of Theorem 4.11

Proof. Constant stepsize

First, by the L_* spectral norm Lipschitz smoothness of f , we have

$$\begin{aligned} f(W_{t+1}) &\leq f(W_t) + \langle \nabla f(W_t), W_{t+1} - W_t \rangle + \frac{L_*}{2} \|W_{t+1} - W_t\|_{\text{op}}^2 \\ &= f(W_t) - \eta_t \|\nabla f(W_t)\|_* + \frac{L_* \eta_t^2}{2} \end{aligned} \quad (13)$$

Then, from the star-convex condition, we obtain

$$\begin{aligned} f(W_t) &\leq f^* + \langle \nabla f(W_t), W_t - W^* \rangle \\ &\leq f^* + \|\nabla f(W_t)\|_* \|W_t - W^*\|_{\text{op}} \\ &\leq f^* + \|\nabla f(W_t)\|_* D_{\text{op}} \\ -\|\nabla f(W_t)\|_* &\leq -\frac{f(W_t) - f^*}{D_{\text{op}}} \end{aligned} \quad (14)$$

where we apply the Cauchy-Schwarz inequality in the second inequality, and in the third inequality we use Assumption 4.10.

Set $\eta = \eta$ and combine (13) and (14). We have

$$f(W_{t+1}) - f^* \leq \left(1 - \frac{\eta}{D_{\text{op}}}\right) (f(W_t) - f^*) + \frac{\eta^2 L_*}{2},$$

and

$$f(W_T) - f^* \leq \left(1 - \frac{\eta}{D_{\text{op}}}\right)^T (f(W_0) - f^*) + \sum_{t=0}^{T-1} \left(1 - \frac{\eta}{D_{\text{op}}}\right)^{T-1-t} \frac{\eta^2 L_*}{2}$$

$$\leq \left(1 - \frac{\eta}{D_{\text{op}}}\right)^T (f(W_0) - f^*) + \frac{\eta L_* D_{\text{op}}}{2}.$$

Thus, to reach the precision $f(W_T) - f^* \leq \epsilon$, we can take $\eta = O(\frac{\epsilon}{L_* D_{\text{op}}})$, and the complexity is $O(L_* D_{\text{op}}^2 \epsilon^{-1} \log \frac{\Delta}{\epsilon})$, where $\Delta = f(W_0) - f^*$.

Adaptive stepsize

By the L_* spectral norm Lipschitz smoothness of f and setting $\eta_t = \frac{\|\nabla f(W_t)\|_*}{L_*}$, we have

$$\begin{aligned} f(W_{t+1}) &\leq f(W_t) + \langle \nabla f(W_t), W_{t+1} - W_t \rangle + \frac{L_*}{2} \|W_{t+1} - W_t\|_{\text{op}}^2 \\ &= f(W_t) - \eta_t \|\nabla f(W_t)\|_* + \frac{L_* \eta_t^2}{2} \\ &= f(W_t) - \frac{\|\nabla f(W_t)\|_*^2}{2L_*}. \end{aligned} \tag{15}$$

From the star-convex property, we have

$$\Delta_t := f(W_t) - f^* \leq \langle \nabla f(W_t), W_t - W^* \rangle \leq \|\nabla f(W_t)\|_* \|W_t - W^*\|_{\text{op}}. \tag{16}$$

Combining Equation (15) and Equation (16), we obtain

$$\Delta_{t+1} \leq \Delta_t - \frac{1}{2L_* \|W_t - W^*\|_{\text{op}}^2} \Delta_t^2.$$

Then we have

$$\frac{1}{\Delta_{t+1}} \geq \frac{1}{\Delta_t} + \frac{\Delta_t}{2L_* \|W_t - W^*\|_{\text{op}}^2 \Delta_{t+1}}.$$

Using $\|W_t - W^*\|_{\text{op}} \leq D_{\text{op}}$, then we can sum above inequality to obtain

$$\frac{1}{\Delta_t} \geq \frac{1}{\Delta_0} + \sum_{i=1}^t \frac{\Delta_{i-1}}{2L_* D_{\text{op}}^2 \Delta_i} \geq \frac{1}{\Delta_0} + \frac{t}{2L_* D_{\text{op}}^2},$$

where we use $\frac{\Delta_t}{\Delta_{t-1}} \geq 1$.

After rearranging terms, we obtain

$$f(W_t) - f^* \leq \frac{2L_*(f(W_0) - f^*)D_{\text{op}}^2}{2L_* D_{\text{op}}^2 + t(f(W_0) - f^*)}.$$

□

C.2 Convergence analysis with Assumption 4.7

Theorem C.1. Under Assumption 4.9, Assumption 4.7 and Assumption 4.10, if we apply Algorithm 2 with $\eta_t = \eta$, we have

$$f(W_T) - f^* \leq \left(1 - \frac{\eta}{D_{\text{op}}}\right)^T (f(W_0) - f^*) + \sum_{t=0}^{T-1} \left(1 - \frac{\eta}{D_{\text{op}}}\right)^{T-1-t} \frac{\eta^2 J_t}{2} + \frac{s\eta^2 D_{\text{op}} r^{3/2}}{6}.$$

where $J_t = \text{vec}(U_t V_t^\top)^\top H_t \text{vec}(U_t V_t^\top)$, $H_t = \nabla^2 f_v(\text{vec}(W_t))$, and $f_v(\text{vec}(W_t)) = f(W_t)$.

Denote $\Delta = f(W_0) - f^*$ and $\hat{J} = \frac{1}{T} \sum_{t=0}^{T-1} |J_t|$. We have

$$f(W_T) - f^* \leq \left(1 - \frac{\eta}{D_{\text{op}}}\right)^T \Delta + \frac{\eta^2 \hat{J} T}{2} + \frac{s\eta^2 D_{\text{op}} r^{3/2}}{6}.$$

When $\hat{J} > 0$, we can set $\eta = \min \left\{ \frac{D_{\text{op}}}{T} \log \left(\frac{T\Delta}{D_{\text{op}}^2 \hat{J}} \right), D_{\text{op}} \right\}$. We have

$$f(W_T) - f^* \leq \frac{D_{\text{op}}^2 \hat{J}}{T} + \frac{D_{\text{op}}^2 \hat{J}}{2T} \left[\log \left(\frac{T\Delta}{D_{\text{op}}^2 \hat{J}} \right) \right]^2 + \frac{s r^{3/2} D_{\text{op}}^3}{6T^2} \left[\log \left(\frac{T\Delta}{D_{\text{op}}^2 \hat{J}} \right) \right]^2$$

$$\leq \tilde{O}\left(\frac{D_{\text{op}}^2 \hat{J}}{T}\right).$$

When $\hat{J} = 0$, we can set $\eta = \min\left\{\frac{D_{\text{op}}}{T} \log\left(\frac{T^2 \Delta}{D_{\text{op}}^3 s r^{3/2}}\right), D_{\text{op}}\right\}$. We have

$$\begin{aligned} f(W_T) - f^* &\leq \frac{s r^{3/2} D_{\text{op}}^3}{T^2} + \frac{s r^{3/2} D_{\text{op}}^3}{6T^2} \left[\log\left(\frac{T^2 \Delta}{D_{\text{op}}^3 s r^{3/2}}\right) \right]^2 \\ &\leq \tilde{O}\left(\frac{s r^{3/2} D_{\text{op}}^3}{T^2}\right). \end{aligned}$$

Proof. According to Equation (12), we have

$$f(W_{t+1}) - f(W_t) \leq -\eta \|\nabla f(W_t)\|_* + \frac{\eta^2 J_t}{2} + \frac{s \eta^3 r^{3/2}}{6} \quad (17)$$

Moreover, according to (14), we have

$$-\|\nabla f(W_t)\|_* \leq -\frac{f(W_t) - f^*}{D_{\text{op}}}$$

Combine (17) and (14), we can obtain

$$f(W_{t+1}) - f^* \leq \left(1 - \frac{\eta}{D_{\text{op}}}\right) (f(W_t) - f^*) + \frac{\eta^2 J_t}{2} + \frac{s \eta^3 r^{3/2}}{6}.$$

Then, we have

$$\begin{aligned} f(W_T) - f^* &\leq \left(1 - \frac{\eta}{D_{\text{op}}}\right)^T (f(W_0) - f^*) + \sum_{t=0}^{T-1} \left(1 - \frac{\eta}{D_{\text{op}}}\right)^{T-1-t} \frac{\eta^2 J_t}{2} + \sum_{t=0}^{T-1} \left(1 - \frac{\eta}{D_{\text{op}}}\right)^{T-1-t} \frac{s \eta^3 r^{3/2}}{6} \\ &\leq \left(1 - \frac{\eta}{D_{\text{op}}}\right)^T (f(W_0) - f^*) + \sum_{t=0}^{T-1} \left(1 - \frac{\eta}{D_{\text{op}}}\right)^{T-1-t} \frac{\eta^2 J_t}{2} + \frac{s \eta^2 D_{\text{op}} r^{3/2}}{6}. \end{aligned}$$

Denote $\Delta = f(W_0) - f^*$ and $\tilde{J} = \sum_{t=0}^{T-1} \left(1 - \frac{\eta}{D_{\text{op}}}\right)^{T-1-t} \frac{\eta}{D_{\text{op}}} J_t$. We have

$$f(W_T) - f^* \leq \left(1 - \frac{\eta}{D_{\text{op}}}\right)^T \Delta + \frac{\eta D_{\text{op}} \tilde{J}}{2} + \frac{s \eta^2 D_{\text{op}} r^{3/2}}{6}.$$

Thus, the convergence complexity can be depending on \tilde{J} , which can be viewed as a weighted average of J_t .

In experiments, we found that usually $J_t > 0$. Note that if f is strict convex and if $U_t V_t^\top \neq 0$, then J_t is strictly larger than 0. Thus, to get a more illustrative and clearer result, if we denote $\hat{J} = \frac{1}{T} \sum_{t=0}^{T-1} |J_t|$, then we have

$$f(W_T) - f^* \leq \left(1 - \frac{\eta}{D_{\text{op}}}\right)^T \Delta + \frac{\eta^2 \hat{J} T}{2} + \frac{s \eta^2 D_{\text{op}} r^{3/2}}{6}.$$

When $\hat{J} > 0$, we can set $\eta = \min\left\{\frac{D_{\text{op}}}{T} \log\left(\frac{T \Delta}{D_{\text{op}}^2 \hat{J}}\right), D_{\text{op}}\right\}$. We have

$$\begin{aligned} f(W_T) - f^* &\leq \frac{D_{\text{op}}^2 \hat{J}}{T} + \frac{D_{\text{op}}^2 \hat{J}}{2T} \left[\log\left(\frac{T \Delta}{D_{\text{op}}^2 \hat{J}}\right) \right]^2 + \frac{s r^{3/2} D_{\text{op}}^3}{6T^2} \left[\log\left(\frac{T \Delta}{D_{\text{op}}^2 \hat{J}}\right) \right]^2 \\ &\leq \tilde{O}\left(\frac{D_{\text{op}}^2 \hat{J}}{T}\right). \end{aligned}$$

When $\hat{J} = 0$, we can set $\eta = \min \left\{ \frac{D_{\text{op}}}{T} \log \left(\frac{T^2 \Delta}{D_{\text{op}}^3 s r^{3/2}} \right), D_{\text{op}} \right\}$. We have

$$\begin{aligned} f(W_T) - f^* &\leq \frac{s r^{3/2} D_{\text{op}}^3}{T^2} + \frac{s r^{3/2} D_{\text{op}}^3}{6T^2} \left[\log \left(\frac{T^2 \Delta}{D_{\text{op}}^3 s r^{3/2}} \right) \right]^2 \\ &\leq \tilde{O} \left(\frac{s r^{3/2} D_{\text{op}}^3}{T^2} \right). \end{aligned}$$

□

C.3 Convergence analysis with Assumption 3.1

Theorem C.2. Under Assumption 4.9, Assumption 3.1, and Assumption 4.10, if we apply Algorithm 2 with $\eta_t = \eta$, then

$$f(W_T) - f^* \leq \left(1 - \frac{\eta}{D_{\text{op}}} \right)^T (f(W_0) - f^*) + \frac{r L D_{\text{op}} \eta}{2}.$$

Thus, to reach the precision $f(W_T) - f^* \leq \epsilon$, we can take $\eta = O(\frac{\epsilon}{r L D_{\text{op}}})$, and the complexity is $O(r L D_{\text{op}}^2 \epsilon^{-1} \log \frac{\Delta}{\epsilon})$, where $\Delta = f(W_0) - f^*$.

Moreover, if we apply Algorithm 2 with adaptive learning rates $\eta_t = \frac{\|\nabla f(W_t)\|_*}{r L}$, we have

$$f(W_t) - f^* \leq \frac{2r L (f(W_0) - f^*) D_{\text{op}}^2}{2r L D_{\text{op}}^2 + t(f(W_0) - f^*)}.$$

Thus, to reach the precision $f(W_T) - f^* \leq \epsilon$, the complexity is $O(r L D_{\text{op}}^2 \epsilon^{-1})$.

Proof. Constant stepsize

First, by the L Frobenius norm Lipschitz smoothness of f , we have

$$\begin{aligned} f(W_{t+1}) &\leq f(W_t) + \langle \nabla f(W_t), W_{t+1} - W_t \rangle + \frac{L}{2} \|W_{t+1} - W_t\|_{\text{F}}^2 \\ &\leq f(W_t) - \eta_t \|\nabla f(W_t)\|_* + \frac{r L \eta_t^2}{2} \end{aligned} \tag{18}$$

Then, according to (14), we have

$$-\|\nabla f(W_t)\|_* \leq -\frac{f(W_t) - f^*}{D_{\text{op}}}$$

Set $\eta = \eta$ and combine (18) and (14). We have

$$f(W_{t+1}) - f^* \leq \left(1 - \frac{\eta}{D_{\text{op}}} \right) (f(W_t) - f^*) + \frac{r L \eta^2}{2},$$

and

$$\begin{aligned} f(W_T) - f^* &\leq \left(1 - \frac{\eta}{D_{\text{op}}} \right)^T (f(W_0) - f^*) + \sum_{t=0}^{T-1} \left(1 - \frac{\eta}{D_{\text{op}}} \right)^{T-1-t} \frac{r L \eta^2}{2} \\ &\leq \left(1 - \frac{\eta}{D_{\text{op}}} \right)^T (f(W_0) - f^*) + \frac{r L D_{\text{op}} \eta}{2}. \end{aligned}$$

Thus, to reach the precision $f(W_T) - f^* \leq \epsilon$, we can take $\eta = O(\frac{\epsilon}{r L D_{\text{op}}})$, and the complexity is $O(r L D_{\text{op}}^2 \epsilon^{-1} \log \frac{\Delta}{\epsilon})$, where $\Delta = f(W_0) - f^*$.

Adaptive stepsize

By the L spectral norm Lipschitz smoothness of f and setting $\eta_t = \frac{\|\nabla f(W_t)\|_*}{r L}$, we have

$$f(W_{t+1}) \leq f(W_t) + \langle \nabla f(W_t), W_{t+1} - W_t \rangle + \frac{L}{2} \|W_{t+1} - W_t\|_{\text{F}}^2$$

$$\begin{aligned}
&\leq f(W_t) - \eta_t \|\nabla f(W_t)\|_* + \frac{rL\eta_t^2}{2} \\
&= f(W_t) - \frac{\|\nabla f(W_t)\|_*^2}{2rL}.
\end{aligned} \tag{19}$$

From the star-convex property, we have

$$\Delta_t := f(W_t) - f^* \leq \langle \nabla f(W_t), W_t - W^* \rangle \leq \|\nabla f(W_t)\|_* \|W_t - W^*\|_{\text{op}}. \tag{20}$$

Combining Equation (19) and Equation (20), we obtain

$$\Delta_{t+1} \leq \Delta_t - \frac{1}{2rL\|W_t - W^*\|_{\text{op}}^2} \Delta_t^2.$$

Then we have

$$\frac{1}{\Delta_{t+1}} \geq \frac{1}{\Delta_t} + \frac{\Delta_t}{2rL\|W_t - W^*\|_{\text{op}}^2 \Delta_{t+1}}.$$

Using $\|W_t - W^*\|_{\text{op}} \leq D_{\text{op}}$, we can sum above inequality to obtain

$$\frac{1}{\Delta_t} \geq \frac{1}{\Delta_0} + \sum_{i=1}^t \frac{\Delta_{i-1}}{2rLD_{\text{op}}^2 \Delta_i} \geq \frac{1}{\Delta_0} + \frac{t}{2rLD_{\text{op}}^2},$$

where we use $\frac{\Delta_t}{\Delta_{t-1}} \geq 1$.

After rearranging terms, we obtain

$$f(W_t) - f^* \leq \frac{2rL(f(W_0) - f^*)D_{\text{op}}^2}{2rLD_{\text{op}}^2 + t(f(W_0) - f^*)}.$$

□

D Experimental settings

Experiments of Figure 1 and Figure 2 are conducted on NVIDIA H100 GPUs. Experiments of Figure 3 are conducted on Intel(R) Core(TM) i9-14900HX.

In Figure 1 (a), we randomly select a fixed subset of 100 samples from CIFAR-10 and train a fully connected neural network (Table 1) on this subset, serving as the deterministic setting. The learning rates for GD, Adam, Muon (using Newton-Schulz iterations with momentum), Muon_withoutM (without momentum), and Muon_SVD (Algorithm 1) are tuned from the set $\{1e^{-4}, 5e^{-4}, 1e^{-3}, 5e^{-3}, 1e^{-2}, 5e^{-2}, 1e^{-1}\}$. We conduct experiments with five different random seeds, and report the average results along with one standard deviation in Figure 1 (a). In Figure 1 (b), we train the same neural network using stochastic batches sampled from the entire CIFAR-10 training set, with a batch size of 120. The learning rates for GD, Muon, Muon_withoutM, and Muon_SVD are tuned from $\{1e^{-3}, 5e^{-3}, 1e^{-2}, 5e^{-2}, 1e^{-1}\}$, while the learning rate for Adam is tuned from $\{1e^{-4}, 5e^{-4}, 1e^{-3}, 5e^{-3}\}$. The model is trained for 5 epochs, and after each epoch, we evaluate the loss on the entire CIFAR-10 training set. The results are presented in Figure 1 (b).

In Figure 2, we randomly select a fixed subset of 60 samples from MNIST. We then train a fully connected neural network (Table 2) with three matrix parameters ($W^1 \in \mathbb{R}^{128 \times 784}$, $W^2 \in \mathbb{R}^{64 \times 128}$, $W^3 \in \mathbb{R}^{10 \times 64}$) on this subset using GD and Muon (Algorithm 2). We first tune the learning rates for GD and Muon from the set $\{1e^{-3}, 5e^{-3}, 1e^{-2}, 5e^{-2}, 1e^{-1}, 5e^{-1}, 1\}$, and find that the optimal learning rate for both is 0.1. Using the selected learning rate, we optimize the neural network with GD and Muon (Algorithm 2), respectively, and record the loss, J_t , L_t , $\|\nabla f(W_t)\|_*^2$, and $\|\nabla f(W_t)\|_{\text{F}}^2$ throughout the optimization process. These quantities are computed with respect to W^2 , i.e., $\nabla f = \nabla_{W^2} f(W_t^1, W_t^2, W_t^3)$ and $\nabla^2 f = \nabla_{\text{vec}(W^2)\text{vec}(W^2)}^2 f_v(\text{vec}(W_t^1), \text{vec}(W_t^2), \text{vec}(W_t^3))$, where $f_v(\text{vec}(W_t^1), \text{vec}(W_t^2), \text{vec}(W_t^3)) = f(W_t^1, W_t^2, W_t^3)$. The results are shown in Figure 2.

In Figure 3(a)(b), we consider the quadratic function $f(W) = \frac{1}{2} \text{Tr}((W - W_{\text{opt}})Q(W - W_{\text{opt}})')$, $W \in \mathbb{R}^{15 \times 20}$, $Q \in \mathbb{R}^{15 \times 15}$, where W_{opt} is randomly chosen and Q is chosen with an ill condition number. Then, we apply GD and Muon to optimize this function, both start from the $W_0 = 0$ with

4000 iterations. We choose the optimal constant stepsize $\frac{1}{L}$ for GD and choose the stepsize for Muon such that Muon can converge in 4000 iterations with the best function value. For each iteration of both algorithms, we record the difference in function value from the optimum and the spectral norm error to the optimal point. The results are shown in Figure 3(a) and (b). In Figure 3(c), we compare the key constant in convergence analysis under quadratic function with ill-conditioned Q . We choose $W_{\text{opt}} \sim U(-50, 50)$, which means each element of W_{opt} is i.i.d. random variables drawn from a continuous uniform distribution. Then, we estimate the key constants $D_{\text{op}}^2 L_*$ and $D_{\text{F}}^2 L$, where we choose the initial point as $W_0 = 0$. We calculate the ratios $\frac{D_{\text{F}}^2 L}{D_{\text{op}}^2 L_*}$ over 10^4 random samples. The results are shown in Figure 3(c).

D.1 Model Architectures

We present the architectures of the models used in our experiments in the following tables (Table 1, 2). For every Linear module, the bias term is set as false – so the models contain only matrix parameters.

Table 1: Model Architecture for Figure 1

Layer Type	Matrix Shape ($m \times n$)
Fully Connected + ReLU	512×3072
Fully Connected + ReLU	256×512
Fully Connected + ReLU	64×256
Fully Connected	10×64

Table 2: Model Architecture for Figure 2

Layer Type	Matrix Shape ($m \times n$)
Fully Connected + ReLU	128×784
Fully Connected + ReLU	64×128
Fully Connected	10×64

E Muon with Newton-Schulz iterations

As discussed in Section 3.1, computing SVD is usually expensive in practice, especially when the dimensions are large. Therefore, a common implementation of Muon employs the Newton-Schulz iterations Bernstein and Newhouse [2024], Jordan et al. [2024] to approximate the orthogonalization process (Algorithm 3). In the original paper of Muon [Jordan et al., 2024], they define their Newton-Schulz iterations (line 6 in Algorithm 3) in the following way. First, the input M_t will be normalized as $X_0 = M_t / \|M_t\|_{\text{F}}$. Then, X_k is computed from X_{k-1} according to the following update rule:

$$X_k = aX_{k-1} + b(X_{k-1}X_{k-1}^\top)X_{k-1} + c(X_{k-1}X_{k-1}^\top)^2X_{k-1}.$$

Following the analysis in Jordan et al. [2024] and supposing the SVD of M_t is $U_t S_t V_t^\top$, we have $X_0 = U_t \bar{S}_t V_t^\top$, where $\bar{S}_t = S_t / \|S_t\|_{\text{F}}$. We can note that

$$\begin{aligned} X_1 &= aX_0 + b(X_0X_0^\top)X_0 + c(X_0X_0^\top)^2X_0 \\ &= (aI + bU\bar{S}^2U^\top + cU\bar{S}^4U^\top)U\bar{S}V^\top \\ &= U(a\bar{S} + b\bar{S}^3 + c\bar{S}^5)V^\top. \end{aligned}$$

Define $\varphi(x) = ax + bx^3 + cx^5$. Then after k iterations, the output is $O_t = X_k = U_t \varphi^k(\bar{S}_t) V_t^\top$. In Jordan et al. [2024], they set $k = 5$, $a = 3.4445$, $b = -4.7750$, and $c = 2.0315$. Note that $\bar{S}_t = S_t / \|S_t\|_{\text{F}}$ is normalized; thus, its singular values are in $[0, 1]$. We can plot $\varphi^5(x)$ in $[0, 1]$ and get Figure 4. Note that there is a steep jump around $x = 0$ in Figure 4 and if we define a threshold ε such that for $x \in [\varepsilon, 1]$, $\varphi^5(x) \in [0.65, 1.25]$, we can have that this threshold ε is very small and most function values of $\varphi^5(x)$ lie in $[0.65, 1.25]$. Thus, $O_t = U_t \varphi^k(\bar{S}_t) V_t^\top$ serves as a low-cost yet effective approximation to $U_t V_t^\top$.

Algorithm 3 Muon with Newton–Schulz iterations

- 1: **Input:** Initial weights W_0 , learning rate schedule $\{\eta_t\}$, $\beta \in [0, 1)$, batch size B
 - 2: **for** $t = 0, 1, \dots, T - 1$ **do**
 - 3: Sample batch $\{\xi_{t,i}\}_{i=1}^B$ uniformly
 - 4: $G_t = \frac{1}{B} \sum_{i=1}^B \nabla f(W_t; \xi_{t,i})$ (or $G_t = \nabla f(W_t)$ in the deterministic setting)
 - 5: If $t > 0$, $M_t = \beta M_{t-1} + (1 - \beta)G_t$. If $t = 0$, $M_0 = G_0$.
 - 6: $O_t = \text{NewtonSchulz}(M_t)$
 - 7: $W_{t+1} = W_t - \eta_t O_t$
 - 8: **end for**
-

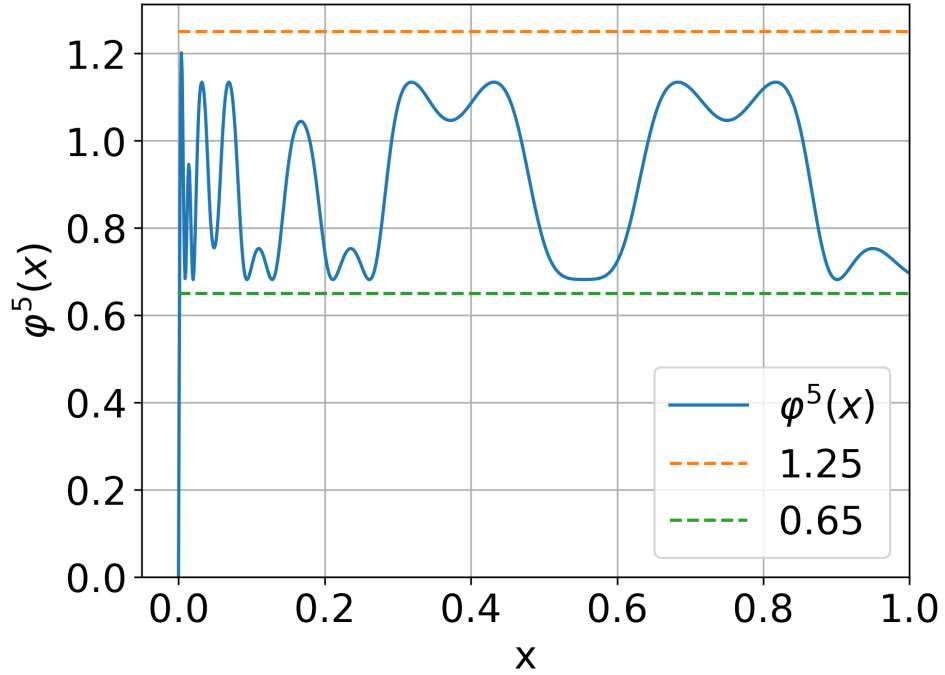


Figure 4: $\varphi^5(x)$ with $\varphi(x) = ax + bx^3 + cx^5$ and $a = 3.4445$, $b = -4.7750$, $c = 2.0315$. Similar to the Figure 4 in [Jordan et al. \[2024\]](#). Line 0.65 and 1.25 are just for illustrative purposes; one can actually choose tighter bounds.