# Leveraging Coordinate Momentum in SignSGD and Muon: Memory-Optimized Zero-Order LLM Fine-Tuning

Grigoriy Evseev
MIPT, Moscow

Egor Petrov
MIPT, Moscow

Veprikov Andrey
MIPT, Moscow
ISP RAS, Moscow

Aleksandr Beznosikov
MIPT, Moscow
Innopolis University, Innopolis

## Abstract

Fine-tuning Large Language Models (LLMs) is essential for adapting pre-trained models to downstream tasks. Yet traditional first-order optimizers such as Stochastic Gradient Descent (SGD) and Adam incur prohibitive memory and computational costs that scale poorly with model size. In this paper, we investigate zero-order (ZO) optimization methods as a memory- and compute-efficient alternative, particularly in the context of parameter-efficient fine-tuning techniques like LoRA. We propose `JAGUAR SignSGD`, a ZO momentum-based algorithm that extends ZO SignSGD, requiring the same number of parameters as the standard ZO SGD and only $\mathcal{O}(1)$ function evaluations per iteration. To the best of our knowledge, this is the first study to establish rigorous convergence guarantees for SignSGD in the stochastic ZO case. We further propose `JAGUAR Muon`, a novel ZO extension of the Muon optimizer that leverages the matrix structure of model parameters, and we provide its convergence rate under arbitrary stochastic noise. Through extensive experiments on challenging LLM fine-tuning benchmarks, we demonstrate that the proposed algorithms meet or exceed the convergence quality of standard first-order methods, achieving significant memory reduction. Our theoretical and empirical results establish new ZO optimization methods as a practical and theoretically grounded approach for resource-constrained LLM adaptation.

## 1 Introduction

Fine-tuning pre-trained Large Language Models (LLMs) has become the standard technique in modern natural language processing [28, 82, 79, 39], enabling rapid adaptation to diverse downstream tasks with minimal labelled data [60, 67, 78]. These models, often trained on massive corpora, achieve state-of-the-art results when fine-tuned on specific applications, including question answering, summarization, and dialogue generation. The fine-tuning setup can be considered as a stochastic unconstrained optimization problem of the form

$$f^* := \min_{x \in \mathbb{R}^d} \left\{ f(x) := \mathbb{E}_{\xi \sim \mathcal{D}} \left[ f(x, \xi) \right] \right\}, \tag{1}$$

where $x$ are parameters of the fine-tuned LLM, $\mathcal{D}$ is the data distribution available for training, and $f(x, \xi)$ is the loss on data point $\xi$.

The de facto standard for solving (1) is the use of First-Order (FO) optimization methods. These approaches assume access to the stochastic gradient $\nabla f(x, \xi)$. Classical FO methods, such as Stochastic Gradient Descent (SGD) [2] and Adam [35], remain the most widely used techniques

for model adaptation due to their efficiency and compatibility with the backpropagation algorithm. Nevertheless, in contemporary fine-tuning tasks, alternative FO algorithms are often preferred.

A recent trend in optimization for LLMs is to represent optimization parameters in matrix form rather than as vectors [6, 5, 58]. Algorithms such as Shampoo [26] and SOAP [73] have demonstrated superior performance on LLM training tasks compared to Adam and SGD [11], which operate in an element-wise manner and do not utilize the underlying structure of the model parameters. Currently, the canonical matrix-based optimization algorithm is Muon [33, 44, 41], which integrates the principles of Shampoo and SOAP but does not employ any preconditioning matrices [33]. The central idea of this method is to project the gradient at each iteration onto the space of semi-orthogonal matrices using the Newton–Schultz algorithm [6].

However, as LLMs continue to scale, the backpropagation procedure, necessary for FO methods, becomes increasingly expensive in terms of memory consumption. For instance, the memory cost of computing gradients during the training of OPT-13B is reported to be more than an order of magnitude larger than that of inference [84]. This imbalance poses a serious bottleneck for deploying LLM fine-tuning in resource-constrained environments such as edge devices [83, 19], consumer-grade GPUs [43, 76], or large-scale distributed settings [27]. To overcome these limitations, researchers are exploring various approaches to reduce the size of the required optimizer statistics. One such approach is the SignSGD algorithm, initially developed for distributed optimization [75], but which has also proven effective in LLM fine-tuning [57], owing to its simplicity, memory efficiency, and surprising empirical effectiveness across a range of adaptation tasks [32, 52]. SignSGD was first rigorously analyzed in the FO setting by [7] and [3]. Minimal memory usage and straightforward hyperparameter tuning make SignSGD an attractive choice for memory-constrained fine-tuning of LLMs ($\sim 4/3\times$ memory usage compared to Adam). Beyond SignSGD, other FO methods also target memory reduction. AdaFactor [69] was among the first, lowering memory usage by storing a single value per block ($\sim 4/3\times$). Additional techniques include quantizing optimizer states to lower-precision formats [13, 40] ($\sim 4/3\times$ and $\sim 16/9\times$ respectively) and fusing the backward pass with optimizer updates [48] ($\sim 4/3\times$), further decreasing memory demands during training.

Nevertheless, the most memory-efficient methods are based on the Zero-Order (ZO) optimization technique, which avoids backpropagation entirely by estimating gradients using only forward passes. This flexibility allows us to treat the model as a black box, optimizing performance with minimal assumptions about its architecture or implementation details. Recent studies [50] have demonstrated the practical benefits of this approach: for example, the MeZO algorithm applies classical ZO SGD [20] to fine-tune LLMs while maintaining four times lower memory requirements than traditional FO methods [51] ($\sim 10\times$ compared to Adam [81]). In ZO methods it is assumed that we only have access to the values of the stochastic function $f(x, \xi)$ from (1) [18, 20]. Within LLMs pretraining or fine-tuning context, oracles are forward passes with small perturbations in parameters of the model. To estimate gradients, authors use finite differences:

$$\nabla f(x, \xi) \approx \widetilde{\nabla} f(x, \xi) = \frac{f(x + \tau e, \xi) - f(x - \tau e, \xi)}{2\tau} e, \tag{2}$$

where $\tau > 0$ is a small number, frequently referred to as a smoothing parameter, and $e \in \mathbb{R}^d$ is some random vector [55, 14, 51, 81]. In the next section, we provide review about different ZO optimization methods, that somehow utilize formula (2).

## 2   Related Work and Our Contributions

**ZO gradient estimators.** The simplest zero-order gradient estimator employs the estimate (2) as the stochastic gradient. However, even this approach presents specific challenges, particularly regarding the selection of an appropriate distribution from which to sample the random vector $e$. The most commonly employed distributions include a uniform sampling over the unit sphere: $e \sim RS(1)_{\|\cdot\|}^d$ [18, 55], a Gaussian distribution with zero mean and identity covariance matrix: $e \sim \mathcal{N}(0, I)$ [55, 20], and standard basis one-hot vectors [14, 68]. Also, some papers [42, 66, 1] utilize the so-called full coordinate estimate, which approximates the gradient across all basis vectors. However, this approach requires $\mathcal{O}(d)$ calls to the zero-order oracle, making it impractical for large-scale fine-tuning tasks. Despite the prevalence of these approaches, alternative and more complicated sampling strategies have also been explored.

2

In [64, 56], the authors explore low-dimensional perturbations within random subspaces. The central concept of random subspace methods involves generating the perturbation vector $e$ within a subspace spanned by a projection matrix $P \in \mathbb{R}^{d \times r}$ and a low-dimensional random vector $\tilde{e} \in \mathbb{R}^r$: $e = P\tilde{e}$. Typically, $P$ and $\tilde{e}$ are sampled from a Gaussian distribution and $r \ll d$. The primary motivation for this method lies in the fact that gradients during the fine-tuning process exhibit a low-dimensional structure [56]. In [47, 74], the authors employ a masked random vector $e$, wherein at each iteration a random mask with $r$ non-zero elements $m_r \in \{0,1\}^d$ is generated and applied element-wise to a Gaussian vector $e$. This procedure accelerates the optimization step, as only the parameters corresponding to the active entries in $m_r$ are updated, rather than the entire parameter set. In contrast, the authors of [25] depart from random mask sampling at each iteration and instead select an optimal mask $m_r$ prior to training, according to a specific criterion. Consequently, the update rule (2) modifies only the parameters selected by the optimal mask during optimization. In our approach, we similarly do not utilize all coordinates of the random vector $e$ in each estimation of (2), instead, we select a single coordinate at each step. However, unlike previous works [47, 74, 25], we do not discard information from the remaining coordinates, but accumulate information from previous iterations. We employ the JAGUAR zero-order gradient estimation technique [72, 53], which integrates the concept of sampling one-hot basis vectors with the utilization of a SAGA-like momentum update [12]. This approach facilitates convergence in the stochastic setting by leveraging memory from past iterations, while using the same amount of memory as standard zero-order methods like ZO SGD (MeZO) [51]. In the original paper [72], the authors do not incorporate a momentum parameter, discarding coordinate information from previous iterations. In contrast, we introduce a momentum parameter, $0 \leq \beta \leq 1$ (see Algorithms 1 and 2), which controls the utilization of gradients from past iterations. We demonstrate that adding this momentum $\beta$ allows the method to converge in the stochastic non-convex case (see Theorems 3.5 and 3.7).

**Momentum techniques.** Numerous zero-order methods in the literature incorporate momentum techniques in various forms. However, these approaches typically introduce multiple additional variables of dimension $d$. Since zero-order methods are often chosen for fine-tuning tasks to save memory, the inclusion of such extra variables becomes a critical limitation in these settings. In [30], authors use variance reduction technique SPIDER [17], that uses approximately $5d$ parameters: $2d$ for ZO gradients, $2d$ for model parameters and $1d$ for momentum. In [9, 31], the authors employ the Adam optimization technique [35], which is frequently used for stochastic non-convex optimization problems [9, 16]. However, this technique incurs a significant memory overhead, requiring $4d$ parameters. The paper [62] utilizes classical heavy-ball momentum within a zero-order framework, provided, only demonstrating almost sure convergence to a constant in the non-convex setting. In our work, we successfully incorporated a momentum technique using only $2d + 1$ parameters and proved the convergence rate within the standard stochastic non-convex setting (see Algorithm 1 and Theorem 3.5). It is worth noting that numerous other zero-order techniques exist in the literature to achieve convergence when the function $f$ is convex [23, 55, 14], satisfies conditions like PL [62] or ABG [61], or in deterministic settings [23]. Since our focus is on fine-tuning problems, which fall under the stochastic non-convex case, we will not discuss these methods in detail.

**Matrix ZO optimization.** In the context of zero-order optimization, transitioning to matrix-valued parameters necessitates replacing the random vector $e \in \mathbb{R}^d$ in zero-order gradient approximation (2) with a random matrix $E \in \mathbb{R}^{m \times n}$, and correspondingly, projecting this matrix $E$ onto a semi-orthogonal space, as is done in the Muon algorithm [33]. Since the random matrix $E$ is typically drawn from a known distribution, it is possible to directly sample orthogonal matrices when computing the gradient estimator (2). A similar approach has previously appeared in the zero-order optimization literature [10]; however, that work did not consider the Muon algorithm, but rather focused on sampling two orthogonal matrices $V \in \mathbb{R}^{m \times r}$ and $U \in \mathbb{R}^{n \times r}$ of rank $r \ll \min\{m, n\}$. This approach does not correspond to the decomposition of the random matrix $E$, as $E$ is almost surely of full rank. Additionally, alternative techniques for sampling low-rank matrices have been proposed in the literature. For instance, in [77], a method analogous to the sampling of low-rank vectors described in [64, 56] is utilized. In our work, we extend our memory-efficient momentum method to the ZO version of the matrix-based Muon algorithm [33] (see Algorithm 2 and Theorem 3.7), keeping the $2d + 1$ parameter efficiency while also broadening our analysis to more modern algorithms that leverage the matrix structure of parameters.

We present a summary of relevant results from the existing zero-order literature in Table 1.

Table 1: Summary of relevant results from the existing zero-order literature.

| | Method | Parameter Count | Convergence Rate Stochastic Non-convex Case | Momentum | Fine-tuning (LLM) Setup |
|---|---|---|---|---|---|
| **Vector Parameters** $\mathbf{x} \in \mathbb{R}^d$ | ZO-SGD [20] | $2 \cdot \mathbf{d}$ | ✓ | ✗ | ✗ |
| | ZO-PSGD [22] | $2 \cdot \mathbf{d}$ | ✓ | ✗ | ✗ |
| | ZO-SCD [42] [(1)] | $2 \cdot \mathbf{d}$ | ✓ | ✗ | ✗[(2)] |
| | ZO-SPIDER [17] | $5 \cdot \mathbf{d}$ | ✓ | ✓ | ✗ |
| | ZO-AdaMM [9] | $4 \cdot \mathbf{d}$ | ✓ | ✓ | ✗ |
| | ZO-SignSGD [45] | $2 \cdot \mathbf{d}$ | ✗✓[(3)] | ✗ | ✗[(4)] |
| | Acc-ZOM [30] | $5 \cdot \mathbf{d}$ | ✓ | ✓ | ✗ |
| | DSFBSD [64] | $(1 + \mathbf{r}) \cdot \mathbf{d}$ [(5)] | ✗ | ✗ | ✗ |
| | MeZO [51] | $2 \cdot \mathbf{d}$ | ✗ | ✗ | ✓ |
| | ZO-ProxSTORM [59] | $5 \cdot \mathbf{d}$ | ✓ | ✓ | ✗ |
| | HB ZO-SGD [62] | $3 \cdot \mathbf{d}$ | ✗[(6)] | ✓ | ✗ |
| | Sparse ZO-SGD [24] | $(2 + \mathbf{r}) \cdot \mathbf{d}$ [(5)] | ✗ | ✗ | ✓ |
| | Sparse MeZO [47] | $3 \cdot \mathbf{d}$ | ✗ | ✗ | ✓ |
| | LeZO [74] | $2 \cdot \mathbf{d}$ | ✗ | ✗ | ✓ |
| | ZO-AdaMU [31] | $4 \cdot \mathbf{d}$ | ✓ | ✓ | ✓ |
| | ZO-SGD-Cons [34] | $2 \cdot \mathbf{d}$ | ✗ | ✗ | ✓ |
| | SGFM [56] | $(2 + \mathbf{r}) \cdot \mathbf{d}$ [(5)] | ✗ | ✗ | ✗ |
| | CompSGD [36] | $2 \cdot \mathbf{d}$ | ✗✓[(3)] | ✗ | ✓ |
| | **JAGUAR SignSGD** **Algorithm** 1 | $2 \cdot \mathbf{d} + 1$ | ✓ | ✓ | ✓ |
| **Matrix Parameters** $\mathbf{X} \in \mathbb{R}^{m \times n}$ | ZO-RMS [49] [(7)] | $2 \cdot \mathbf{mn}$ | ✗✓[(3)] | ✗ | ✗ |
| | MeZO [51] | $2 \cdot \mathbf{mn}$ | ✗ | ✗ | ✓ |
| | LOZO [10] | $(\mathbf{m + n})\mathbf{r} + 2 \cdot \mathbf{mn}$ [(5)] | ✓ | ✗ | ✓ |
| | SubZero [77] [(8)] | $(\mathbf{m + n + r})\mathbf{r} + 2 \cdot \mathbf{mn}$ [(5)] | ✗ | ✗ | ✓ |
| | **JAGUAR Muon** **Algorithm** 2 | $2 \cdot \mathbf{mn} + 1$ | ✓ | ✓ | ✓ |

[(1)] Uses a full coordinate ZO estimator. [(2)] Considers asynchronous algorithms. [(3)] Convergence only to a neighborhood of the solution. [(4)] Addresses adversarial attacks in deep learning. [(5)] $r \ll d, m, n$ is a small number. [(6)] Only asymptotic convergence to a constant. [(7)] Assumes that parameters are symmetric matrices. [(8)] Assumes sparsity of parameters.

## 2.1 Our Contributions

While zero-order optimization methods have recently attracted attention for LLM fine-tuning, previous work has primarily focused on basic algorithms. In this paper, we broaden the scope of zero-order optimization by introducing advanced momentum techniques, specifically adapting the JAGUAR approach [72] to the SignSGD algorithm in the zero-order setting (see Algorithms 1). We consider this algorithm because SignSGD has demonstrated state-of-the-art performance in LLM fine-tuning tasks, outperforming even AdamW [57]. Our key contributions are as follows:

- We provide the first convergence analysis in the stochastic non-convex setting for zero-order SignSGD with momentum (Algorithm 1 and Theorem 3.5), requiring only $2d + 1$ parameters and $\mathcal{O}(1)$ ZO oracle calls per iteration.

- We extend our memory-efficient momentum method to the Muon algorithm (Algorithm 2), introducing the first zero-order variant of Muon that preserves memory efficiency. We also establish its convergence rate in the stochastic non-convex setting (Theorem 3.7).

- We empirically evaluate the proposed zero-order methods on challenging LLM fine-tuning benchmarks, demonstrating their effectiveness and practical relevance.

# 3 Main results

## 3.1 Preliminaries

**Notations.** We denote the $\ell_1$ and $\ell_2$ (Euclidean) norms of a vector $x \in \mathbb{R}^d$ as $\|x\|_1 := \sum_{i=1}^{d} |x_i|$ and $\|x\|_2^2 := \sum_{i=1}^{d} x_i^2$, respectively. For clarity, matrix-valued variables are denoted by capital letters.

For matrices $X \in \mathbb{R}^{m \times n}$, we use the Schatten 1-norm ($\mathcal{S}_1$) and Schatten 2-norm ($\mathcal{S}_2$, Frobenius): $\|X\|_{\mathcal{S}_1} := \sum_{i=1}^{d} |(\Sigma_X)_{i,i}|$ and $\|X\|_{\mathcal{S}_2}^2 := \sum_{i=1}^{d} (\Sigma_X)_{i,i}^2 = \sum_{i=1}^{m} \sum_{j=1}^{n} X_{i,j}^2 =: \|X\|_F^2$, where $X = U_X \Sigma_X V_X^T$ is the reduced Singular Value Decomposition (SVD) of $X$. The standard dot product between two vectors $x, y \in \mathbb{R}^d$ is defined as $\langle x, y \rangle := x^T y$. For matrices $X, Y \in \mathbb{R}^{m \times n}$, we define the inner product as $\langle X, Y \rangle := \text{tr}(X^T Y)$.

We now provide several assumptions that are necessary for the analysis.

**Assumption 3.1** (Smoothness). The functions $f(x, \xi)$ are $L(\xi)$-smooth on the $\mathbb{R}^d$ with respect to the Euclidean norm $\|\cdot\|$, i.e., for all $x, y \in \mathbb{R}^d$ it holds that
$$\|\nabla f(x, \xi) - \nabla f(y, \xi)\|_2 \le L(\xi)\|x - y\|_2.$$
We also assume that exists constant $L^2 := \mathbb{E}\left[L(\xi)^2\right]$.

**Assumption 3.2** (Bounded variance of the gradient). The variance of the $\nabla f(x, \xi)$ is bounded with respect to the Euclidean norm, i.e., there exists $\sigma > 0$, such that for all $x \in \mathbb{R}^d$ it holds that
$$\mathbb{E}\left[\|\nabla f(x, \xi) - \nabla f(x)\|_2^2\right] \le \sigma^2.$$

We assume access only to a zero-order oracle, which returns a noisy evaluation of the function $f(x, \xi)$. Therefore, we are limited to using this noisy value $\hat{f}(x, \xi)$ in the estimation of the ZO gradient (2). This noise may originate not only from inherent randomness (stochastic noise), but also from systematic effects (deterministic noise), such as computer rounding errors. Therefore, we make a common assumption about the function $\hat{f}(x, \xi)$ returned by the oracle [15, 72].

**Assumption 3.3** (Bounded oracle noise). The noise in the oracle is bounded with respect to the Euclidean norm, i.e., there exists $\Delta > 0$, such that for all $x \in \mathbb{R}^d$ it holds that
$$\mathbb{E}\left[\left|\hat{f}(x, \xi) - f(x, \xi)\right|^2\right] \le \Delta^2.$$

Assumptions 3.1 and 3.2 are standard in the theoretical analysis of stochastic non-convex zero-order optimization problems [62, 25, 47, 74]. In contrast, Assumption 3.3 is frequently omitted in the existing literature, as it is commonly presumed that $\Delta = 0$, implying access to an ideal zero-order oracle. However, this assumption does not hold in practice, as numerical errors such as machine precision inevitably introduce a non-zero perturbation. Consequently, while $\Delta$ is typically small, it is never zero, which does not allow us to restore a true gradient along the direction $e$ in the estimation (2) if we set $\tau \to 0$.

## 3.2 Zero-Order Momentum SignSGD with JAGUAR Gradient Approximation

In this section, we introduce zero-order SignSGD algorithm with JAGUAR gradient approximation [72, 53] and momentum of the form:

---

**Algorithm 1** Zero-Order Momentum SignSGD with JAGUAR (`JAGUAR SignSGD`)

---

1: **Parameters:** stepsize (learning rate) $\gamma$, momentum $\beta$, gradient approximation parameter $\tau$, number of iterations $T$.
2: **Initialization:** choose $x^0 \in \mathbb{R}^d$ and $m^{-1} = \mathbf{0} \in \mathbb{R}^d$.
3: **for** $t = 0, 1, 2, \ldots, T$ **do**
4:     Sample $i_t \sim \text{Uniform}(\overline{1, d})$
5:     Set one-hot vector $e^t$ with 1 in the $i_t$ coordinate: $e_{i_t}^t = 1$ and $e_{i \ne i_t}^t = 0$ for all $i \in \overline{1, d}$
6:     Sample stochastic variable $\xi^t \sim \mathcal{D}$
7:     Compute $\widetilde{\nabla}_{i_t} f(x^t, \xi^t) := \frac{\hat{f}(x^t + \tau e^t, \xi^t) - \hat{f}(x^t - \tau e^t, \xi^t)}{2\tau} \in \mathbb{R}$
8:     Set $m_{i_t}^t = \beta m_{i_t}^{t-1} + (1 - \beta)\widetilde{\nabla}_{i_t} f(x^t, \xi^t)$ and $m_{i \ne i_t}^t = m_{i \ne i_t}^{t-1}$ for all $i \in \overline{1, d}$
9:     Set $x^{t+1} = x^t - \gamma \cdot \text{sign}(m^t)$
10: **end for**
11: **Return:** $x^{N(T)}$, where $N(T) \sim \text{Uniform}(\overline{1, T})$.

---

The gradient approximation employed in Algorithm 1 deviates from that of the original JAGUAR method, as we introduce a momentum variable $\beta$. The estimator from the original work can be recovered by setting $\beta = 0$.

We now present a lemma characterizing the closeness between the momentum variable $m^t$ from line 8 and the true gradient $\nabla f(x^t)$.

**Lemma 3.4.** *Consider $m^t$ from line 8 of Algorithm 1. Under Assumptions 3.1, 3.2, 3.3 it holds that:*

$$\mathbb{E}\left[\left\|m^t - \nabla f(x^t)\right\|_2^2\right] = \mathcal{O}\left[\frac{d^3 L^2 \gamma^2}{(1-\beta)^2} + (1-\beta)d\sigma^2 + dL^2\tau^2 + \frac{2d\Delta^2}{\tau^2} + \left(\frac{1-\beta}{d}\right)^t \left\|\nabla f(x^0)\right\|_2^2\right].$$

**Discussion.** This lemma closely parallels Lemma 1 from [72], with the key distinction that our analysis incorporates the momentum parameter $\beta$, which was not present in [72]. The introduction of momentum is essential for proving convergence of algorithms such as SignSGD (Algorithm 1) and Muon (see Algorithm 2 in the next section) in the stochastic zero-order setting [71], as it enables more careful handling of variance $\sigma$ in the gradient estimates (2). Another important difference from prior works is that result from Lemma 3.4 does not involve the term $\|\nabla f(x^t)\|_2^2$, which typically appears in analyses where the zero-order gradient estimator (2) is constructed using random uniform or Gaussian vectors $e$ [8, 38, 23, 59]. With the presence of the term $\|\nabla f(x^t)\|_2^2$, it is not possible to achieve convergence guarantees for SignSGD (Algorithm 1) and Muon (Algorithm 2) even with momentum in the stochastic zero-order setting. It is worth noting that a similar result can be obtained when using a full coordinate estimator [42]. However, this approach requires $\mathcal{O}(d)$ calls to the zero-order oracle per iteration, which can be computationally expensive. In contrast, the JAGUAR method achieves the same result with only $\mathcal{O}(1)$ oracle calls and with the same number of parameters, offering significant improvements in efficiency. This makes our approach particularly attractive for large-scale optimization tasks, where reducing oracle complexity is critical.

With the help of Lemma 3.4, we provide convergence analysis of `JAGUAR SignSGD` (Algorithm 1).

**Theorem 3.5.** *Consider Assumptions 3.1, 3.2 and 3.3. Then `JAGUAR SignSGD` (Algorithm 1) has the following convergence rate:*

$$\mathbb{E}\left[\left\|\nabla f\left(x^{N(T)}\right)\right\|_1\right] = \mathcal{O}\left[\frac{\delta_0}{\gamma T} + \frac{d\left\|\nabla f(x^0)\right\|_2}{T\sqrt{1-\beta}} + \frac{d^2 L\gamma}{1-\beta} + \sqrt{1-\beta}d\sigma + dL\tau + \frac{d\Delta}{\tau}\right],$$

*where we used a notation $\delta_0 := f(x^0) - f^*$.*

**Corollary 3.6.** *Consider the conditions of Theorem 3.5. In order to achieve the $\varepsilon$-approximate solution (in terms of $\mathbb{E}\left[\left\|\nabla f(x^{N(T)})\right\|_1\right] \leq \varepsilon$), Algorithm 1 needs $T$ iterations (ZO oracle calls), for:*

***Arbitrary tuning:*** $\gamma = \gamma_0 \cdot T^{-3/4}d^{-1}$, $\beta = 1 - T^{-1/2}$, $\tau = (\Delta/L)^{1/2}$ *and* $\varepsilon \geq d\sqrt{\Delta L}$ :

$$T = \mathcal{O}\left[\left(\frac{d\delta_0/\gamma_0 + d\left\|\nabla f(x^0)\right\|_2 + dL\gamma_0 + d\sigma}{\varepsilon}\right)^4\right].$$

***Optimal tuning:*** $\gamma = \sqrt{\frac{\delta_0(1-\beta)}{d^2 LT}}$, $\beta = 1 - \min\left\{1; \sqrt{\frac{L\delta_0}{T\sigma^2}}\right\}$, $\tau = (\Delta/L)^{1/2}$ *and* $\varepsilon \geq d\sqrt{\Delta L}$ :

$$T = \mathcal{O}\left[\frac{\delta_0 L d^2}{\varepsilon^2} + \frac{\delta_0 L d^2}{\varepsilon^2} \cdot \left(\frac{d\sigma}{\varepsilon}\right)^2\right].$$

**Discussion.** The convergence rate established in Theorem 3.5 is similar to what is known for first-order methods [7, 32, 65, 36], however our bounds include an additional factor of $d$, which is typical for all coordinate-based methods [54, 63], not just zero-order ones. This dependence on the dimension arises because coordinate methods process one direction at a time, accumulating complexity proportional to $d$. It is also important to note that without momentum ($\beta = 0$), the algorithm can only guarantee convergence to a neighbourhood of the optimum of size proportional to $\sigma$, as shown in previous works on zero-order SignSGD [45, 36]. Let us also point out that we cannot choose an arbitrary $\varepsilon$ in Corollary 3.6, since there exists an irreducible [15, 72] error $\Delta$ in the zero-order oracle (see Assumption 3.3). However, since $\Delta$ is very small, we can still achieve an acceptable accuracy $\varepsilon$. In our analysis, we use the $\ell_1$-norm of the gradient as the convergence criterion, while the standard in non-convex optimization is the $\ell_2$-norm (Euclidean) [20, 21]. By

setting $\varepsilon_{\ell_1} = \sqrt{d} \cdot \varepsilon_{\ell_2}$, we can rescale our result of Corollary 3.6 for optimal tuning (one can easily do a similar transformation for arbitrary tuning) as

$$T_{\text{Euclidean}} = \mathcal{O}\left[\frac{\delta_0 L d}{\varepsilon^2} + \frac{\delta_0 L d}{\varepsilon^2} \cdot \left(\frac{\sqrt{d}\sigma}{\varepsilon}\right)^2\right].$$

This substitution allows us to obtain improved results in terms of the dependence on $d$.

### 3.3 Zero-Order Muon with JAGUAR Gradient Approximation

In this section, we address the matrix optimization setting, where the optimization variables $X_t$ are elements of the matrix space $\mathbb{R}^{m \times n}$, rather than the standard vector space $\mathbb{R}^d$. Such a formulation allows for a more direct representation of model parameters, helping to better capture their underlying structure [6, 58]. For the first time in the literature, we introduce a zero-order version of the Muon [33] algorithm (Algorithm 2), broadening the applicability to matrix-structured optimization tasks where only function evaluations are available.

---

**Algorithm 2** Zero-Order Muon with JAGUAR (`JAGUAR Muon`)

---

1: **Parameters:** stepsize (learning rate) $\gamma$, momentum $\beta$, gradient approximation parameter $\tau$, number of Newton-Schulz steps ns_steps, number of iterations $T$.
2: **Initialization:** choose $X^0 \in \mathbb{R}^{m \times n}$ and $M^{-1} = \mathbf{0} \in \mathbb{R}^{m \times n}$.
3: **for** $t = 0, 1, 2, \ldots, T$ **do**
4:      Sample $i_t \sim \text{Uniform}(\overline{1, m})$ and $j_t \sim \text{Uniform}(\overline{1, n})$
5:      Set one-hot matrix $E^t$ with 1 in the $(i_t, j_t)$ coordinate
6:      Sample stochastic variable $\xi^t \sim \mathcal{D}$
7:      Compute $\widetilde{\nabla}_{i_t} f(X^t, \xi^t) := \frac{\hat{f}(X^t + \tau E^t, \xi^t) - \hat{f}(X^t - \tau E^t, \xi^t)}{2\tau} \in \mathbb{R}$
8:      Set $M^t_{i_t, j_t} = \beta M^{t-1}_{i_t, j_t} + (1-\beta)\widetilde{\nabla}_{i_t} f(x^t, \xi^t)$ and $M^t_{i \neq i_t, j \neq j_t} = M^{t-1}_{i \neq i_t, j \neq j_t}$
9:      Set $X^{t+1} = X^t - \gamma \cdot \texttt{Newton\_Schulz}(M^t, K = \text{ns\_steps})$
10: **end for**
11: **Return:** $X^{N(T)}$, where $N(T) \sim \text{Uniform}(\overline{1, T})$.
1: **Subroutine** $\texttt{Newton\_Schulz}(A \in \mathbb{R}^{m \times n}, K = 10)$ [6]:
2:      Set $A^0 = A/\|A\|_F$
3:      **for** $k = 0, 1, 2, \ldots, K$ **do**
4:          $A^{k+1} = 3/2 \cdot A^k - 1/2 \cdot A^k (A^k)^T A^k$
5:      **end for**
6:      **Return:** $A^K \approx U_A \cdot V_A^T$.

---

Algorithm 2 is similar to the first-order Muon algorithm [33], the only difference is that we use zero-order gradient approximation JAGUAR [72] in line 8.

Let us note that when extending to matrix-valued parameters, it is necessary to slightly modify Assumptions 3.1 and 3.2: all occurrences of the $\ell_2$ norm $\|\cdot\|_2$ should be replaced with the Frobenius norm $\|\cdot\|_F$. This modification is justified, as the following property holds for all matrices $A \in \mathbb{R}^{m \times n}$: $\|A\|_F = \|\overline{\text{vec}}(A)\|_2$. We now provide the convergence analysis of `JAGUAR Muon` (Algorithm 2).

**Theorem 3.7.** *Consider Assumptions 3.1, 3.2 (with Frobenius norm) and 3.3. Then* `JAGUAR Muon` *(Algorithm 2) has the following convergence rate:*

$$\mathbb{E}\left[\left\|\nabla f\left(X^{N(T)}\right)\right\|_{\mathcal{S}_1}\right] = \mathcal{O}\left[\frac{\delta_0}{\gamma T} + \frac{m^{1/2} n \left\|\nabla f(X^0)\right\|_2}{T\sqrt{1-\beta}} + \frac{m^{3/2} n^2 \gamma}{1-\beta} + \sqrt{1-\beta} m^{1/2} n \sigma \right.$$
$$\left. + m^{1/2} n L \tau + \frac{m^{1/2} n \Delta}{\tau}\right],$$

*where we used a notation $\delta_0 := f(x^0) - f^*$. We also assume that $n \leq m$.*

**Corollary 3.8.** *Consider the conditions of Theorem 3.7. In order to achieve the $\varepsilon$-approximate solution (in terms of $\mathbb{E}[\|\nabla f(X^{N(T)})\|_{\mathcal{S}_1} \leq \varepsilon$), Algorithm 2 needs $T$ iterations (ZO calls), for:*

**Arbitrary tuning:** $\gamma = \gamma_0 \cdot T^{-3/4}(mn)^{-1}, \beta = 1 - T^{-1/2}, \tau = (\Delta/L)^{1/2}, \varepsilon \geq m^{1/2}n\sqrt{\Delta L}:$

$$T = \mathcal{O}\left[\left(\frac{mn\delta_0/\gamma_0 + m^{1/2}n\left\|\nabla f(X^0)\right\|_2 + m^{1/2}nL\gamma_0 + m^{1/2}n\sigma}{\varepsilon}\right)^4\right].$$

**Optimal tuning:** $\gamma = \sqrt{\frac{\delta_0(1-\beta)}{m^{3/2}n^2 LT}}, \beta = 1 - \min\left\{1; \sqrt{\frac{L\delta_0}{T\sigma^2}}\right\}, \tau = (\Delta/L)^{1/2}, \varepsilon \geq m^{1/2}n\sqrt{\Delta L}:$

$$T = \mathcal{O}\left[\frac{\delta_0 L m^{3/2}n^2}{\varepsilon^2} + \frac{\delta_0 L m^{3/2}n^2}{\varepsilon^2} \cdot \left(\frac{m^{3/2}n^2\sigma}{\varepsilon}\right)^2\right].$$

**Discussion.** The convergence rate established in Theorem 3.7 is consistent with the first-order case [41, 37]. However, there remain zero-order terms depending on $\tau$ and $\Delta$, as for Algorithm 1 (see Theorem 3.5 and Discussion part after it). From a proof perspective, Theorems 3.5 and 3.7 are very similar, since the orthogonalization operation (`Newton_Schulz`) in Algorithm 2 can be interpreted as taking the sign of the gradient matrix eigenvalues. Accordingly, both the form and the convergence rate criterion are analogous (the $\ell_1$ norm for Algorithm 1 and the $\mathcal{S}_1$ norm for Algorithm 2). Nevertheless, the convergence rates of the two algorithms differ slightly. We examine the two boundary cases in the following remark.

*Remark* 3.9. For optimal tuning from Corollary 3.8 we can specify the number of iterations of Algorithm 2 to achieve the $\varepsilon$-approximate solution in terms of the total number of parameters $d = m \cdot n$ in the two boundary cases:

- If $n \ll m \approx d$:

$$T_{n \ll m \approx d} = \mathcal{O}\left[\frac{\delta_0 L d^{3/2}}{\varepsilon^2} + \frac{\delta_0 L d^{3/2}}{\varepsilon^2} \cdot \left(\frac{d^{3/2}\sigma}{\varepsilon}\right)^2\right].$$

- If $n \approx m \approx \sqrt{d}$:

$$T_{n \approx m \approx \sqrt{d}} = \mathcal{O}\left[\frac{\delta_0 L d^{7/4}}{\varepsilon^2} + \frac{\delta_0 L d^{7/4}}{\varepsilon^2} \cdot \left(\frac{d^{7/4}\sigma}{\varepsilon}\right)^2\right].$$

Accordingly, comparing these convergence rates with that obtained in Corollary 3.8, we observe an improvement by factors of $d^{1/2}$ and $d^{1/4}$, respectively.

# 4 Experiments

In this section, we present a comprehensive empirical evaluation to validate the theoretical contributions of our proposed ZO optimization methods for fine-tuning large language models. Our study aims to assess both the accuracy and memory efficiency of these methods, comparing them against established ZO and FO baselines. We build upon the experimental framework proposed in [81], extending it to incorporate our novel algorithms: `JAGUAR SignSGD` (Algorithm 1) and `JAGUAR Muon` (Algorithm 2). The primary objective is to achieve competitive test accuracy on downstream tasks while maintaining memory efficiency comparable to the baseline methods. Additionally, we introduce `ZO-Muon` (Algorithm 3 in Appendix A), a direct zero-order adaptation of Muon [33], utilizing the standard Gaussian zero-order gradient estimation (2). We also explore various matrix sampling techniques to enhance this approach, as detailed in Appendix A.

## 4.1 Experimental Setup

**Fine-Tuning Task and Schemes.** Fine-tuning LLMs is a pivotal process in adapting pre-trained models to specific downstream tasks, enabling high performance with limited task-specific data. This process typically involves adjusting model parameters to minimize a task-specific loss function, a task that becomes computationally intensive as model sizes scale to billions of parameters. To explore the efficacy of our ZO methods, we focus on the SST2 dataset [70], a widely-used benchmark for binary sentiment classification [81, 10, 51]. We evaluate two fine-tuning schemes:

- **Full Fine-Tuning (FT):** Updates all parameters of the pre-trained model, offering maximum flexibility at the cost of higher computational resources.

- **Low-Rank Adaptation (LoRA):** Introduces a small set of trainable parameters while keeping the original model parameters frozen, enhancing memory efficiency [29].

**Models.** We conduct experiments using three prominent LLMs: OPT-1.3B [80], a 1.3 billion parameter model from the OPT family; RoBERTa-Large [46], a 355 million parameter model known for its robust performance in natural language processing tasks. These models represent a range of sizes and architectures, allowing us to assess the scalability and generality of our methods.

**Methods.** We evaluate the following ZO optimization methods proposed in this work:

- `JAGUAR SignSGD`: Combines the JAGUAR gradient approximation [72] with SignSGD and momentum for efficient updates (Algorithm 1).

- `JAGUAR Muon`: Integrates JAGUAR with the Muon optimizer, incorporating momentum and orthogonalization (Algorithm 2).

- `ZO-Muon`: A novel ZO adaptation of the Muon optimizer, leveraging matrix-based optimization principles (Algorithm 3 in Appendix A).

**Comparison procedure.** For comparison, we include baseline methods from [81]: ZO-SGD [20], Acc-ZOM [30], ZO-SGD-Cons [34], ZO-SignSGD [45], ZO-AdaMM [9], Forward-Grad [4], and the FO method FO-SGD [2]. The results for which are given in the benchmark paper. Additionally, we compare our methods with LeZO [74], which employs a comparable layer-wise selection mechanism similar to `JAGUAR SignSGD` coordinate-wise updates. We perform experiments for our methods in accordance with similar experiments from [81]. For details of our hyperparameter selection and model training procedure, see Appendix.

## 4.2 Results

**Accuracy Comparison.** Table 2 presents the test accuracy results for SST2 across both models and fine-tuning schemes. Our proposed methods demonstrate strong performance, often outperforming baseline ZO methods.

**Discussion.** For OPT-1.3B with LoRA, ZO Jaguar achieves the highest accuracy of $94.0\% \pm 0.1$, surpassing ZO-AdaMM's 92.3%, while `JAGUAR SignSGD` and `JAGUAR Muon` reach $92.5\% \pm 0.5$ and $93.5\% \pm 0.1$, respectively. In the FT setting, `JAGUAR SignSGD` excels with $94.0\% \pm 0.1$, significantly improving over ZO-SGD's 90.8%. However, `ZO-Muon` and `JAGUAR Muon` exhibit lower FT performance at $86.5\% \pm 0.1$ and $72.5\% \pm 0.1$, respectively, potentially due to the presence of non-matrix parameters in the full FT process. For RoBERTa-Large, `JAGUAR SignSGD` achieves $92.2\% \pm 0.2$ in FT and $92.2\% \pm 0.4$ in LoRA, consistently outperforming most baselines, `JAGUAR Muon` reached the same $92.2\% \pm 0.2$ in LoRA, while `ZO-Muon` maintains competitive results.

Table 2: Test accuracy on SST2 for OPT-1.3B and RoBERTa-Large with FT and LoRA. Best performance among ZO methods is in **bold**. Blue indicates outperformance of all baseline ZO methods, red indicates matching or exceeding FO-SGD.

| Method | OPT-1.3B | | RoBERTa-Large | |
|---|---|---|---|---|
| | FT | LoRA | FT | LoRA |
| FO-SGD | 91.1 | 93.6 | 91.4 | 91.2 |
| Forward-Grad | 90.3 | 90.3 | 90.1 | 89.7 |
| ZO-SGD | 90.8 | 90.1 | 89.4 | 90.8 |
| Acc-ZOM | 85.2 | 91.3 | 89.6 | 90.9 |
| ZO-SGD-Cons | 88.3 | 90.5 | 89.6 | 91.6 |
| ZO-SignSGD | 87.2 | 91.5 | 52.5 | 90.2 |
| ZO-AdaMM | 84.4 | 92.3 | 89.8 | 89.5 |
| LeZO | 85.1 | 92.3 | 90.4 | 91.8 |
| JAGUAR SignSGD | **94.0 ± 0.1** | 92.5 ± 0.5 | **92.2 ± 0.2** | 92.2 ± 0.4 |
| ZO-Muon | 86.5 ± 0.1 | 93.5 ± 0.1 | 85.0 ± 0.1 | 91.0 ± 0.2 |
| JAGUAR Muon | 72.5 ± 0.1 | **94.0 ± 0.1** | 72.0 ± 0.1 | **92.2 ± 0.2** |

## 4.3 Memory Efficiency

Table 3 compares peak memory usage for OPT-1.3B, highlighting the efficiency of our methods. Our methods consume approximately 4.17-4.18 GB in FT and 4.12-4.13 GB in LoRA, closely aligning with ZO-SGD, while FO-SGD requires 12.2 GB and 5.9 GB, respectively. This demonstrates that our approaches effectively balance accuracy gains with memory efficiency.

Table 3: Peak memory usage (GB) for OPT-1.3B (half-precision, F16) on SST2 with FT and LoRA.

| Method | FT Memory | LoRA Memory |
|---|---|---|
| FO-SGD | 12.246 | 5.855 |
| ZO-SGD | 4.171 | 4.125 |
| ZO-AdaMM | 9.046 | 4.132 |
| JAGUAR SignSGD | 4.172 | 4.128 |
| ZO-Muon | 4.177 | 4.130 |
| JAGUAR Muon | 4.179 | 4.132 |

# References

[1] Zeeshan Akhtar and Ketan Rajawat. Zeroth and first order stochastic frank-wolfe algorithms for constrained optimization. *IEEE Transactions on Signal Processing*, 70:2119–2135, 2022.

[2] Shun-ichi Amari. Backpropagation and stochastic gradient descent method. *Neurocomputing*, 5(4):185–196, 1993.

[3] Lukas Balles and Philipp Hennig. Dissecting adam: The sign, magnitude and variance of stochastic gradients. In *International Conference on Machine Learning (ICML)*, 2017.

[4] Atılım Güneş Baydin, Barak A. Pearlmutter, Don Syme, Frank Wood, and Philip Torr. Gradients without backpropagation, 2022.

[5] Jeremy Bernstein and Laker Newhouse. Modular duality in deep learning. *arXiv preprint arXiv:2410.21265*, 2024.

[6] Jeremy Bernstein and Laker Newhouse. Old optimizer, new norm: An anthology. *arXiv preprint arXiv:2409.20325*, 2024.

[7] Jeremy Bernstein, Jiawei Zhao, Kamyar Azizzadenesheli, and Anima Anandkumar. signsgd with majority vote is communication efficient and fault tolerant. *arXiv preprint arXiv:1810.05291*, 2018.

[8] HanQin Cai, Yuchen Lou, Daniel McKenzie, and Wotao Yin. A zeroth-order block coordinate descent algorithm for huge-scale black-box optimization. In *ICML*, pages 1182–1191, 2021.

[9] Xiangyi Chen, Sijia Liu, Kaidi Xu, Xingguo Li, Xue Lin, Mingyi Hong, and David Cox. Zo-adamm: Zeroth-order adaptive momentum method for black-box optimization. *Advances in neural information processing systems*, 32, 2019.

[10] Yiming Chen, Yuan Zhang, Liyuan Cao, Kun Yuan, and Zaiwen Wen. Enhancing zeroth-order fine-tuning for language models with low-rank structures. *ArXiv*, abs/2410.07698, 2024.

[11] George E Dahl, Frank Schneider, Zachary Nado, Naman Agarwal, Chandramouli Shama Sastry, Philipp Hennig, Sourabh Medapati, Runa Eschenhagen, Priya Kasimbeg, Daniel Suo, et al. Benchmarking neural network training algorithms. *arXiv preprint arXiv:2306.07179*, 2023.

[12] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. *Advances in neural information processing systems*, 27, 2014.

[13] Tim Dettmers, Mike Lewis, Sam Shleifer, and Luke Zettlemoyer. 8-bit optimizers via block-wise quantization. *arXiv preprint arXiv:2110.02861*, 2021.

[14] John C. Duchi, Michael I. Jordan, Martin J. Wainwright, and Andre Wibisono. Optimal rates for zero-order convex optimization: The power of two function evaluations. *IEEE Transactions on Information Theory*, 61(5):2788–2806, 2015.

[15] Pavel Dvurechensky, Eduard Gorbunov, and Alexander Gasnikov. An accelerated directional derivative method for smooth stochastic convex optimization. *European Journal of Operational Research*, 290(2):601–621, 2021.

[16] Anonymous Author et al. Mezo-a$^3$dam: Memory-efficient zeroth-order adam with adaptivity adjustments. *OpenReview, ICLR 2025*, 2024. Under review.

[17] Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang. Spider: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. *Advances in neural information processing systems*, 31, 2018.

[18] Abraham D. Flaxman, Adam Tauman Kalai, and H. Brendan McMahan. Online convex optimization in the bandit setting: gradient descent without a gradient. In *Proceedings of the Sixteenth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 385–394, 2005.

[19] Lei Gao, Amir Ziashahabi, Yue Niu, Salman Avestimehr, and Murali Annavaram. Enabling efficient on-device fine-tuning of llms using only inference engines. *arXiv preprint arXiv:2409.15520*, 2024.

[20] Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM journal on optimization*, 23(4):2341–2368, 2013.

[21] Saeed Ghadimi and Guanghui Lan. Accelerated gradient methods for nonconvex optimization. *arXiv preprint arXiv:1608.06860*, 2016.

[22] Saeed Ghadimi, Guanghui Lan, and Hongchao Zhang. Mini-batch stochastic approximation methods for nonconvex stochastic composite optimization. *Mathematical Programming*, 155(1):267–305, 2016.

[23] Eduard Gorbunov, Pavel Dvurechensky, and Alexander Gasnikov. An accelerated method for derivative-free smooth stochastic convex optimization. *SIAM Journal on Optimization*, 32(2):1210–1238, 2022.

[24] Wentao Guo, Jikai Long, Yimeng Zeng, Zirui Liu, Xinyu Yang, Yide Ran, Jacob R Gardner, Osbert Bastani, Christopher De Sa, Xiaodong Yu, et al. Zeroth-order fine-tuning of llms with extreme sparsity. *arXiv preprint arXiv:2406.02913*, 2024.

[25] Wentao Guo, Jikai Long, Yimeng Zeng, Zirui Liu, Xinyu Yang, Yide Ran, Jacob R. Gardner, Osbert Bastani, Christopher De Sa, Xiaodong Yu, Beidi Chen, and Zhaozhuo Xu. Zeroth-order fine-tuning of llms with extreme sparsity, 2024.

[26] Vineet Gupta, Tomer Koren, and Yoram Singer. Shampoo: Preconditioned stochastic tensor optimization, 2018.

[27] Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint*, abs/1510.00149, 2015.

[28] Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. *arXiv preprint arXiv:1801.06146*, 2018.

[29] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.

[30] Feihu Huang, Shangqian Gao, Jian Pei, and Heng Huang. Accelerated zeroth-order and first-order momentum methods from mini to minimax optimization. *Journal of Machine Learning Research*, 23(36):1–70, 2022.

[31] Shuoran Jiang, Qingcai Chen, Youcheng Pan, Yang Xiang, Yukang Lin, Xiangping Wu, Chuanyi Liu, and Xiaobao Song. Zo-adamu optimizer: Adapting perturbation by the momentum and uncertainty in zeroth-order optimization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 18363–18371, 2024.

[32] Richeng Jin, Yufan Huang, Xiaofan He, Huaiyu Dai, and Tianfu Wu. Stochastic-sign sgd for federated learning with theoretical guarantees. *arXiv preprint arXiv:2002.10940*, 2020.

[33] 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.

[34] Bumsu Kim, Daniel McKenzie, HanQin Cai, and Wotao Yin. Curvature-aware derivative-free optimization. *Journal of Scientific Computing*, 103(2), March 2025.

[35] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[36] Nikita Kornilov, Philip Zmushko, Andrei Semenov, Alexander Gasnikov, and Alexander Beznosikov. Sign operator for coping with heavy-tailed noise: High probability convergence bounds with extensions to distributed optimization and comparison oracle. *arXiv preprint arXiv:2502.07923*, 2025.

[37] Dmitry Kovalev. Understanding gradient orthogonalization for deep learning via non-euclidean trust-region optimization. *arXiv preprint arXiv:2503.12645*, 2025.

[38] David Kozak, Cesare Molinari, Lorenzo Rosasco, Luis Tenorio, and Silvia Villa. Zeroth order optimization with orthogonal random directions. *arXiv preprint*, abs/2107.03941, 2021.

[39] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021.

[40] Bingrui Li, Jianfei Chen, and Jun Zhu. Memory efficient optimizers with 4-bit states. *Advances in Neural Information Processing Systems*, 36:15136–15171, 2023.

[41] Jiaxiang Li and Mingyi Hong. A note on the convergence of muon and further, 2025.

[42] Xiangru Lian, Huan Zhang, Cho-Jui Hsieh, Yijun Huang, and Ji Liu. A comprehensive linear speedup analysis for asynchronous stochastic parallel optimization from zeroth-order to first-order. *Advances in neural information processing systems*, 29, 2016.

[43] Changyue Liao, Mo Sun, Zihan Yang, Jun Xie, Kaiqi Chen, Binhang Yuan, Fei Wu, and Zeke Wang. Lohan: Low-cost high-performance framework to fine-tune 100b model on a consumer gpu. *arXiv preprint arXiv:2403.06504*, 2024.

[44] Jingyuan Liu, Jianlin Su, Xingcheng Yao, Zhejun Jiang, Guokun Lai, Yulun Du, Yidao Qin, Weixin Xu, Enzhe Lu, Junjie Yan, Yanru Chen, Huabin Zheng, Yibo Liu, Shaowei Liu, Bohong Yin, Weiran He, Han Zhu, Yuzhi Wang, Jianzhou Wang, Mengnan Dong, Zheng Zhang, Yongsheng Kang, Hao Zhang, Xinran Xu, Yutao Zhang, Yuxin Wu, Xinyu Zhou, and Zhilin Yang. Muon is scalable for llm training, 2025.

[45] Sijia Liu, Pin-Yu Chen, Xiangyi Chen, and Mingyi Hong. signsgd via zeroth-order oracle. In *International conference on learning representations*, 2019.

[46] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach, 2019.

[47] Yong Liu, Zirui Zhu, Chaoyu Gong, Minhao Cheng, Cho-Jui Hsieh, and Yang You. Sparse mezo: Less parameters for better performance in zeroth-order llm fine-tuning. *arXiv preprint arXiv:2402.15751*, 2024.

[48] Kai Lv, Yuqing Yang, Tengxiao Liu, Qinghui Gao, Qipeng Guo, and Xipeng Qiu. Full parameter fine-tuning for large language models with limited resources. *arXiv preprint arXiv:2306.09782*, 2023.

[49] Alejandro I Maass, Chris Manzie, Iman Shames, and Hayato Nakada. Zeroth-order optimization on subsets of symmetric matrices with application to mpc tuning. *IEEE Transactions on Control Systems Technology*, 30(4):1654–1667, 2021.

[50] Bharath Malladi, Amir Aghazadeh, Haotian Tang, et al. Mezo: Memory-efficient zeroth-order optimization of large language models. *arXiv preprint*, abs/2305.18660, 2023.

[51] Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D Lee, Danqi Chen, and Sanjeev Arora. Fine-tuning language models with just forward passes. *Advances in Neural Information Processing Systems*, 36:53038–53075, 2023.

[52] Emanuele Mengoli, Luzius Moll, Virgilio Strozzi, and El-Mahdi El-Mhamdi. On the byzantine fault tolerance of signsgd with majority vote. *arXiv preprint arXiv:2502.19170*, 2025.

[53] Ruslan Nazykov, Aleksandr Shestakov, Vladimir Solodkin, Aleksandr Beznosikov, Gauthier Gidel, and Alexander Gasnikov. Stochastic frank-wolfe: Unified analysis and zoo of special cases. In *International Conference on Artificial Intelligence and Statistics*, pages 4870–4878. PMLR, 2024.

[54] Yurii Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.

[55] Yurii Nesterov and Vladimir Spokoiny. Random gradient-free minimization of convex functions. *Foundations of Computational Mathematics*, 17(2):527–566, 2017.

[56] Ryota Nozawa, Pierre-Louis Poirion, and Akiko Takeda. Zeroth-order random subspace algorithm for non-smooth convex optimization. *Journal of Optimization Theory and Applications*, 204(3):53, 2025.

[57] Hanyang Peng, Shuang Qin, Fangqing Jiang, Yue Yu, Hui Wang, and Ge Li. Softsignsgd (s3): An enhanced optimize for practical dnn training and loss spikes minimization beyond adam.

[58] Thomas Pethick, Wanyun Xie, Kimon Antonakopoulos, Zhenyu Zhu, Antonio Silveti-Falls, and Volkan Cevher. Training deep learning models with norm-constrained lmos. *arXiv preprint arXiv:2502.07529*, 2025.

[59] Yuxiang Qian and Yong Zhao. Zeroth-order proximal stochastic recursive momentum algorithm for nonconvex nonsmooth optimization. In *2023 International Conference on New Trends in Computational Intelligence (NTCI)*, volume 1, pages 419–423. IEEE, 2023.

[60] 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.

[61] Marco Rando, Cesare Molinari, Silvia Villa, and Lorenzo Rosasco. Stochastic zeroth order descent with structured directions. *Computational Optimization and Applications*, pages 1–37, 2024.

[62] Tadipatri Uday Kiran Reddy and Mathukumalli Vidyasagar. Convergence of momentum-based heavy ball method with batch updating and/or approximate gradients. In *2023 Ninth Indian Control Conference (ICC)*, pages 182–187. IEEE, 2023.

[63] Peter Richtárik and Martin Takáč. Distributed coordinate descent method for learning with big data. *Journal of Machine Learning Research*, 17(75):1–25, 2016.

[64] Lindon Roberts and Clément W Royer. Direct search based on probabilistic descent in reduced spaces. *SIAM Journal on Optimization*, 33(4):3057–3082, 2023.

[65] Mher Safaryan and Peter Richtárik. Stochastic sign descent methods: New algorithms and better theory. In *International Conference on Machine Learning*, pages 9224–9234. PMLR, 2021.

[66] Anit Kumar Sahu, Manzil Zaheer, and Soummya Kar. Towards gradient free and projection free stochastic optimization. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 3468–3477. PMLR, 2019.

[67] Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. Multitask prompted training enables zero-shot task generalization. *arXiv preprint arXiv:2110.08207*, 2021.

[68] Ohad Shamir. On the complexity of bandit and derivative-free stochastic convex optimization. In *ICML*, pages 1001–1009, 2013.

[69] Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, pages 4596–4604. PMLR, 2018.

[70] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Y. Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, 2013.

[71] Tao Sun, Qingsong Wang, Dongsheng Li, and Bao Wang. Momentum ensures convergence of signsgd under weaker assumptions. In *International Conference on Machine Learning*, pages 33077–33099. PMLR, 2023.

[72] Andrey Veprikov, Alexander Bogdanov, Vladislav Minashkin, and Aleksandr Beznosikov. New aspects of black box conditional gradient: Variance reduction and one point feedback. *Chaos, Solitons & Fractals*, 189:115654, 2024.

[73] 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.

[74] Fei Wang, Li Shen, Liang Ding, Chao Xue, Ye Liu, and Changxing Ding. Simultaneous computation and memory efficient zeroth-order optimizer for fine-tuning large language models. *arXiv preprint arXiv:2410.09823*, 2024.

[75] Haibo Yang, Xin Zhang, Minghong Fang, and Jia Liu. Adaptive multi-hierarchical signsgd for communication-efficient distributed optimization. In *2020 IEEE 21st International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, pages 1–5. IEEE, 2020.

[76] Junjie Yin, Jiahao Dong, Yingheng Wang, Christopher De Sa, and Volodymyr Kuleshov. Modulora: finetuning 2-bit llms on consumer gpus by integrating with modular quantizers. *arXiv preprint arXiv:2309.16119*, 2023.

[77] Ziming Yu, Pan Zhou, Sike Wang, Jia Li, and Hua Huang. Zeroth-order fine-tuning of llms in random subspaces, 2024.

[78] Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. *arXiv preprint arXiv:2106.10199*, 2021.

[79] Hongyi Zhang, Zuchao Li, Ping Wang, and Hai Zhao. Selective prefix tuning for pre-trained language models. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 2806–2813, 2024.

[80] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, et al. Opt: Open pre-trained transformer language models, 2022.

[81] Yihua Zhang, Pingzhi Li, Junyuan Hong, Jiaxiang Li, Yimeng Zhang, Wenqing Zheng, Pin-Yu Chen, Jason D. Lee, Wotao Yin, Mingyi Hong, Zhangyang Wang, Sijia Liu, and Tianlong Chen. Revisiting zeroth-order optimization for memory-efficient llm fine-tuning: A benchmark, 2024.

[82] Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. Dialogpt: Large-scale generative pre-training for conversational response generation. *arXiv preprint arXiv:1911.00536*, 2019.

[83] Ligeng Zhu, Lanxiang Hu, Ji Lin, Wei-Ming Chen, Wei-Chen Wang, Chuang Gan, and Song Han. Pockengine: Sparse and efficient fine-tuning in a pocket. In *Proceedings of the 56th Annual IEEE/ACM International Symposium on Microarchitecture*, pages 1381–1394, 2023.

[84] Yujia Zhu, Yujing Zhang, Ziwei Zhang, et al. Efficient fine-tuning of language models via zeroth-order optimization. *arXiv preprint*, abs/2305.14395, 2023.

# A  Classical ZO Muon

## A.1  Zero-Order Muon

The modification presented here reduces memory usage by eliminating the need to store or compute gradients, instead enabling efficient parameter updates through directional finite differences. By substituting the gradient with its zero-order estimate $G_t = \frac{\mathcal{L}(W_t + \tau E) - \mathcal{L}(W_t - \tau E)}{2\tau} E$, where $E$ is sampled from the standard Normal distribution $\mathcal{N}_{st}$ (each value $E_{i,j} \sim \mathcal{N}(0, 1)$), we adapt the Muon algorithm [33] into its zero-order form:

---

**Algorithm 3** Zero-Order Muon (`ZO-Muon`)

---

1: **Parameters:** stepsize (learning rate) $\gamma$, gradient approximation parameter $\tau$, number of iterations $T$.
2: **Initialization:** choose $W_0 \in \mathbb{R}^{n \times d}$
3: **for** $t = 0, 1, 2, \ldots, T$ **do**
4:     Sample $E \in \mathbb{R}^{n \times d}$ from $\mathcal{N}_{st}$
5:     Compute $\mathbf{G}_t = \frac{\mathcal{L}(W_t + \tau E) - \mathcal{L}(W_t - \tau E)}{2\tau} E$
6:     Compute $\mathbf{O}_t = \text{NewtonSchulz}(\mathbf{G}_t)$
7:     Set $W_{t+1} = W_t - \gamma \mathbf{O}_t$
8: **end for**
1: **Subroutine** `Newton_Schulz`$(A \in \mathbb{R}^{m \times n}, K = 10)$ [6]:
2:     Set $A^0 = A/\|A\|_F$
3:     **for** $k = 0, 1, 2, \ldots, K$ **do**
4:         $A^{k+1} = 3/2 \cdot A^k - 1/2 \cdot A^k (A^k)^T A^k$
5:     **end for**
6:     **Return:** $A^K \approx U_A \cdot V_A^T$.

---

## A.2  Matrix-Efficient Zero-Order Muon

In this section, we consider a zero-order gradient estimate:

$$G = \frac{\mathcal{L}(W + \tau E) - \mathcal{L}(W - \tau E)}{2\tau} E, \tag{3}$$

where $E \in \mathbb{R}^{n \times d}$ is typically sampled from $RS_2(1)$.

Using the Newton-Schulz process [33], we approximate the gradient estimate with its SVD components, where $G = U\Sigma V^T$ yields:

$$\text{NewtonSchulz}(G) \approx UV^T. \tag{4}$$

From (3), we express $E$ as:

$$E = \text{sign}(\mathcal{L}(W + \tau E) - \mathcal{L}(W - \tau E))U \left| \frac{2\tau}{\mathcal{L}(W + \tau E) - \mathcal{L}(W - \tau E)} \right| \Sigma V^T. \tag{5}$$

Given the SVD of $E = U_E \Sigma_E V_E^T$, equating with (5) gives:

$$U = \text{sign}(\mathcal{L}(W + \tau E) - \mathcal{L}(W - \tau E))U_E,$$
$$\Sigma = \left| \frac{2\tau}{\mathcal{L}(W + \tau E) - \mathcal{L}(W - \tau E)} \right| \Sigma_E,$$
$$V = V_E. \tag{6}$$

Substituting into (4), we obtain:

$$\text{NewtonSchulz}(G) \approx \text{sign}(\mathcal{L}(W + \tau E) - \mathcal{L}(W - \tau E))U_E V_E^T. \tag{7}$$

This approximation allows orthogonalization of $G$ using the SVD of $E$ without directly applying Newton-Schulz.

The resulting Matrix-Efficient Zero-Order Muon algorithm is:

---
**Algorithm 4** Matrix-Efficient Zero-Order Muon
---
 1: **Parameters:** stepsize (learning rate) $\gamma$, gradient approximation parameter $\tau$, number of iterations $T$.
 2: **Initialization:** choose $W^0 \in \mathbb{R}^{n \times d}$
 3: **for** $t = 0, 1, 2, \ldots, T$ **do**
 4:     Sample orthogonal $U_E \in \mathbb{R}^{n \times n}$, $V_E \in \mathbb{R}^{d \times d}$, $\Sigma_E \sim \mathcal{D} \in \mathbb{R}^{n \times d}$
 5:     Set $E = U_E \Sigma_E V_E^T$
 6:     Compute $\mathbf{O}_t = \text{sign}(\mathcal{L}(W^t + \tau E) - \mathcal{L}(W^t - \tau E)) U_E V_E^T$
 7:     Set $W^{t+1} = W^t - \gamma \cdot \mathbf{O}_t$
 8: **end for**
---

## A.3 Evaluation Procedure

**Hyperparameter Tuning.** To ensure optimal performance, we conducted a grid search over key hyperparameters for each method:

- Momentum parameter: $\beta \in \{10^{-3}, 10^{-2}\}$,
- Learning rate: $\gamma \in [10^{-6}, 10^{-1}]$,
- Smoothing parameter: $\tau \in \{10^{-1}, 10^{-2}, 10^{-3}\}$.

Additional fixed parameters include a momentum of 0.9 and an epsilon of $10^{-3}$ for numerical stability. The best-performing hyperparameters for each algorithm are detailed in the supplementary material, with notable settings including $\tau = 10^{-3}$, $\beta = 10^{-2}$ across all methods, and method-specific learning rates (e.g., $\gamma = 10^{-3}$ for JAGUAR SignSGD in LoRA).

**Evaluation Metrics.** We assess performance using:

- **Test Accuracy:** Measured as the percentage of correct predictions on the SST2 test set, reflecting model effectiveness.
- **Peak Memory Usage:** Quantified in gigabytes (GB) during training, indicating memory efficiency.

**Implementation Details.** Experiments were conducted with three independent runs per configuration, each with a randomly selected seed fixed at the start to ensure reproducibility. We report the mean and standard deviation of test accuracy. Following [51], we employed half-precision (F16) training for ZO methods and mixed-precision (FP16) training for FO methods to optimize memory usage. Training was performed on a single NVIDIA A100 GPU and a single NVIDIA H100 GPU, with memory profiling conducted using standard PyTorch utilities.

## A.4 Experimental Methodology

Our experimental procedure was designed to rigorously evaluate the proposed methods under controlled conditions. For each combination of dataset (SST2), model (OPT-1.3B, RoBERTa-Large), fine-tuning scheme (FT, LoRA), and optimization method, we executed the following steps:

1. **Initialization:** Loaded the pre-trained model and initialized trainable parameters (all for FT, LoRA-specific for LoRA).

2. **Hyperparameter Selection:** Performed a preliminary parameter search to identify the best hyperparameters per method, iterating over the specified ranges and selecting based on validation accuracy.

3. **Evaluation:** Computed test accuracy on the SST2 test set after each run, averaging results across three runs with different seeds.

4. **Memory Profiling:** Recorded peak memory usage during training, ensuring consistency by maintaining identical hardware settings.

This methodology ensures a fair comparison across methods, capturing both performance and resource utilization comprehensively.

## B  Proofs for ZO Momentum SignSGD with JAGUAR (Algorithm 1)

### B.1  Proof of Theorem 3.5

*Proof.* We start from using Lemma 1 from [71]. For the points $x^t$, generated by Algorithm 1 it holds that:

$$f(x^{t+1}) - f(x^t) \leq -\gamma \left\| \nabla f(x^t) \right\|_1 + 2\sqrt{d}\gamma \left\| m^t - \nabla f(x^t) \right\|_2 + \frac{dL\gamma^2}{2}. \tag{8}$$

Now we take mathematical expectation of the both sides of the inequality (8) and use the results from Lemma 3.4:

$$\mathbb{E}\left[f(x^{t+1})\right] - \mathbb{E}\left[f(x^t)\right] \leq -\gamma \mathbb{E}\left[\left\| \nabla f(x^t) \right\|_1\right] + 2\sqrt{d}\gamma \mathbb{E}\left[\left\| m^t - \nabla f(x^t) \right\|_2\right] + \frac{dL\gamma^2}{2}$$

$$= -\gamma \mathbb{E}\left[\left\| \nabla f(x^t) \right\|_1\right] + \mathcal{O}\left[\frac{d^2}{1-\beta} \cdot L\gamma^2 + \sqrt{1-\beta}d\gamma\sigma + d\gamma L\tau\right.$$

$$\left. + \frac{d\gamma\Delta}{\tau} + \sqrt{d}\gamma\left(1 - \frac{1-\beta}{d}\right)^{t/2} \left\| \nabla f(x^0) \right\|_2\right] + \frac{dL\gamma^2}{2}.$$

Consequently, after summing all $T$ steps, we obtain:

$$\gamma \sum_{t=0}^{T} \mathbb{E}\left[\left\| \nabla f(x^t) \right\|_1\right] = \mathcal{O}\left[f(x^0) - f(x^T) + T \cdot \left(\frac{d^2}{1-\beta} \cdot L\gamma^2 + \sqrt{1-\beta}d\gamma\sigma + d\gamma L\tau\right)\right.$$

$$\left. + T \cdot \frac{d\gamma\Delta}{\tau} + \sqrt{d}\gamma \sum_{t=0}^{T} \left(1 - \frac{1-\beta}{d}\right)^{t/2} \left\| \nabla f(x^0) \right\|_2\right]. \tag{9}$$

Now, we divide equation (9) by $\gamma T$ from both sides and obtain:

$$\frac{1}{T} \sum_{t=0}^{T} \mathbb{E}\left[\left\| \nabla f(x^t) \right\|_1\right] = \mathcal{O}\left[\frac{\delta_0}{\gamma T} + \frac{d\left\| \nabla f(x^0) \right\|_2}{T\sqrt{1-\beta}} + \frac{d^2 L\gamma}{1-\beta} + \sqrt{1-\beta}d\sigma + dL\tau + \frac{d\Delta}{\tau}\right],$$

where we used a notation $\delta_0 := f(x^0) - f^*$. This finishes the proof. $\qquad\square$

## C  Proofs for ZO Muon with JAGUAR (Algorithm 2)

### C.1  Technical Lemmas

**Lemma C.1.** *Consider two arbitrary matrixes $A$, $B$ of the same shape and their SVD decomposition: $A = U_A \Sigma_A V_A^T$, $B = U_B \Sigma_B V_B^T$. Define $r_A$ and $r_B$ as ranks of $A$ and $B$, then it holds that*

$$\left|\left\langle A, U_A V_A^T - U_B V_B^T \right\rangle\right| \leq 2 \left\| A - B \right\|_{\mathcal{S}_1} \leq 2\sqrt{\mathrm{rank}(A - B)} \left\| A - B \right\|_F.$$

*Proof.* We first provide an axillary notation:

$$\delta := \left\langle A, U_A V_A^T - U_B V_B^T \right\rangle.$$

Because $U_A$ and $V_A$ have orthonormal columns:

$$\langle A, U_A V_A^\top \rangle = \mathrm{tr}\left(V_A \Sigma_A U_A^\top U_A V_A^\top\right) = \mathrm{tr}(\Sigma_A) = \|A\|_{\mathcal{S}_1}.$$

Hence

$$\delta = \|A\|_{\mathcal{S}_1} - \langle A, U_B V_B^\top \rangle.$$

Insert $B$ and regroup:

$$\delta = \|A\|_{\mathcal{S}_1} - \left(\langle B, U_B V_B^\top \rangle + \langle A - B, U_B V_B^\top \rangle\right) = \|A\|_{\mathcal{S}_1} - \|B\|_{\mathcal{S}_1} - \langle A - B, U_B V_B^\top \rangle.$$

The first difference is controlled by the triangle inequality for the nuclear norm:

$$\left|\|A\|_{\mathcal{S}_1} - \|B\|_{\mathcal{S}_1}\right| \leq \|A - B\|_{\mathcal{S}_1}.$$

17

618  For the second term, Hölder's inequality with $\|U_B V_B^\top\|_2 = 1$ gives

$$\left|\langle A - B, U_B V_B^\top\rangle\right| \leq \|A - B\|_{\mathcal{S}_1}.$$

619  Therefore

$$|\delta| \leq \|A - B\|_{\mathcal{S}_1} + \|A - B\|_{\mathcal{S}_1} = 2\|A - B\|_{\mathcal{S}_1}.$$

620  Using the connection between the Frobenius ($\mathcal{S}_2$) by nuclear ($\mathcal{S}_1$) norms we obtain that:

$$|\delta| = \left\langle A, U_A V_A^T - U_B V_B^T\right\rangle \leq 2\|A - B\|_{\mathcal{S}_1} \leq 2\sqrt{\mathrm{rank}(A - B)}\|A - B\|_F.$$

621  The factor 2 in the nuclear norm bound is sharp, as equality holds for $B = -A$. This finishes the
622  proof. $\qquad\square$

623  We now provide lemma similar to the step Lemma 1 from [71], but in the matrix case.

624  **Lemma C.2** (Step lemma for Muon with momentum). *Let $f$ be an $L$-smooth function (Assumption*
625  *3.1), and let $X^\dagger, M \in \mathbb{R}^{m \times n}$ with $m \geq n$ be an arbitrary matrixes. We define*

$$X^\ddagger := X^\dagger - \gamma \cdot U_M V_M^T,$$

626  *where $\gamma > 0$ and $U_M V_M^T$ comes from SVD decomposition of $M$: $M = U_M \Sigma_M V_M^T$. Then, it holds*
627  *that:*

$$f\left(X^\ddagger\right) - f\left(X^\dagger\right) \leq -\gamma\left\|\nabla f\left(X^\dagger\right)\right\|_{\mathcal{S}_1} + 2\sqrt{n}\gamma\left\|\nabla f\left(X^\dagger\right) - M\right\|_F + \frac{Ln\gamma^2}{2}.$$

628  *Proof.* The $L$-smoothness of the gradient (Assumption 3.1) gives us

$$\begin{aligned}
f\left(X^\ddagger\right) - f\left(X^\dagger\right) &\leq \left\langle\nabla f\left(X^\dagger\right), X^\ddagger - X^\dagger\right\rangle + \frac{L}{2}\left\|X^\ddagger - X^\dagger\right\|_F^2 \\
&= -\gamma\left\langle\nabla f\left(X^\dagger\right), U_M V_M^T\right\rangle + \frac{Ln\gamma^2}{2} \\
&= -\gamma\left\langle\nabla f\left(X^\dagger\right), U_\nabla V_\nabla^T\right\rangle + \gamma\left\langle\nabla f\left(X^\dagger\right), U_\nabla V_\nabla^T - U_M V_M^T\right\rangle + \frac{Ln\gamma^2}{2},
\end{aligned}$$

629  where $U_\nabla V_\nabla^T$ comes from SVD decomposition of $\nabla f\left(X^\dagger\right)$: $\nabla f\left(X^\dagger\right) = U_\nabla \Sigma_\nabla V_\nabla^T$. Therefore the
630  first dot product takes form:

$$-\gamma\left\langle\nabla f\left(X^\dagger\right), U_\nabla V_\nabla^T\right\rangle = -\gamma\mathrm{tr}\left(V_\nabla \Sigma_\nabla U_\nabla^T U_\nabla V_\nabla^T\right) = -\gamma\mathrm{tr}\left(\Sigma_\nabla\right) = -\gamma\left\|\nabla f\left(X^\dagger\right)\right\|_{\mathcal{S}_1}.$$

631  Now we utilize Lemma C.1 with $A = \nabla f\left(X^\dagger\right)$ and $B = M$:

$$\begin{aligned}
f\left(X^\ddagger\right) - f\left(X^\dagger\right) &\leq -\gamma\left\|\nabla f\left(X^\dagger\right)\right\|_{\mathcal{S}_1} + 2\gamma\left\|\nabla f\left(X^\dagger\right) - M\right\|_{\mathcal{S}_1} + \frac{Ln\gamma^2}{2} \\
&\leq -\gamma\left\|\nabla f\left(X^\dagger\right)\right\|_{\mathcal{S}_1} + 2\sqrt{n}\gamma\left\|\nabla f\left(X^\dagger\right) - M\right\|_F + \frac{Ln\gamma^2}{2}.
\end{aligned}$$

632  This finishes the proof. $\qquad\square$