

THE KY FAN NORMS AND BEYOND: DUAL NORMS AND COMBINATIONS FOR MATRIX OPTIMIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

In this article, we explore the use of matrix norms for optimizing functions of weight matrices, a crucial problem in training large language models. Moving beyond the spectral norm that underlies the Muon update, we leverage the Ky Fan k -norm, its affine combinations with other norms, and their corresponding duals to develop a new family of Muon-like algorithms. We complement our theoretical analysis with an extensive empirical study of the algorithms across a wide range of tasks and settings.

1 INTRODUCTION

Minimizing loss functions in unprecedentedly high-dimensional spaces has recently become an integral and crucial part in training large language models. Hence, new scalable, time- and memory-efficient algorithms have been demanded. Besides well-known Adam and AdamW Kingma & Ba (2014), Loshchilov & Hutter (2017), recently proposed Muon has shown promising results on training very large models Jordan et al. (2024); Liu et al. (2025). Its key difference from Adam and AdamW is that it has been constructed specifically for optimizing functions of weight matrices, which are common in deep learning.

That is what can be said from a practical point of view. From the perspective of theory, Muon’s main innovation was an intentional usage of matrix norms, i.e. the spectral norm, to derive the algorithm’s update Bernstein (2025). Based upon recent theoretical advances that explain some theory behind Muon, Scion and Gluon (Bernstein (2025); Kovalev (2025); Pethick et al. (2025b); Riabinin et al. (2025)), we explore application of other matrix norms to optimization of functions of matrices. As it has been done with Muon, we stipulate that our algorithms’ updates be fast to compute.

In this article, we focus on the two most common matrix norms akin to the spectral, namely, the nuclear Norm and the Frobenius norm. Working in the linear minimization oracle (lmo) framework, which is equivalent to a factor to the trust region approach and the steepest descent under norm constraint, we derive Neon, our algorithm based on the nuclear norm. In the section *Matrix side of the updates*, we explain how Neon updates can be computed asymptotically faster than Muon updates by the Newton-Schulz iterations.

Noticing that Neon and Muon are diametrical in terms of the rank of the update matrix, we bridge the space by “regularizing” them by NormalizedSGD, which is derived in lmo with the Frobenius norm. We do this in the same lmo approach by considering a norm that is dual to the convex combination of the Frobenius norm and the spectral or the nuclear norms respectively. So we derive the algorithms we name F-Neon and F-Muon respectively.

Table that compares methods, updates, and sends to the proofs?

Having faced the array of different Muon-like optimizers, according to the upper bounds from Kovalev (2025); Riabinin et al. (2025), with similar convergence behavior, we painstakingly compare them on a synthetic linear least squares problem with known Lipschitz constant. The efforts results in comparison of the algorithms by their convergence not in terms of the dual norms of their gradient, but in terms of the common spectral norm, which may strongly differ, especially in large matrices, from the initial dual norms. Thus, we compare the algorithms in a unified fashion.

Finally, we test Muon, Neon, NSGD, F-Muon and F-Neon on deep-learning tasks: training convolutional network on CIFAR-10 and fine-tuning NanoGPT. The results support the supremacy of Muon,

but the most striking result of the tests is that F-Muon, only half of which is Muon, surpasses Muon's accuracy on the CIFAR tasks by a margin. The case of F-Muon answers in the affirmative to the question of feasibility of constructing a mixture of optimization algorithms to increase robustness of the composite algorithm.

2 PROBLEM STATEMENT

2.1 FUNCTION OF A MATRIX AND ASSUMPTIONS

Warning: the text was copied, while formulas were adapted. We need to rephrase the subsection! We consider the problem of minimizing a differentiable matrix function $F(\cdot): \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$

$$\min_{\mathbf{X} \in \mathbb{R}^{m \times n}} f(\mathbf{X}). \quad (1)$$

To make the theoretical analysis possible in the future, we make three reasonable non-restrictive assumptions. Idea: move gradient assumptions and smoothness to the L-smooth part. Here we need only custom norm and their relation to the Frobenius norm.

Non-Euclidean norm setting and Lipschitz continuous gradient. We assume that matrix space $\mathbb{R}^{m \times n}$ is equipped with a norm $\|\cdot\|: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}_+$, which possibly does not coincide with the Frobenius, norm $\|\cdot\|_F$. In addition, we assume that the gradient $\nabla f(\cdot)$ is Lipschitz continuous with respect to the norm $\|\cdot\|$, that is, the following inequality holds:

$$\|\nabla f(\mathbf{X}) - \nabla f(\mathbf{X}')\|^\dagger \leq L \|\mathbf{X} - \mathbf{X}'\| \quad \text{for all } \mathbf{X}, \mathbf{X}' \in \mathbb{R}^{m \times n}, \quad (A1)$$

where $L > 0$ is the gradient Lipschitz constant, and $\|\cdot\|^\dagger: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}_+$ is the dual norm associated with $\|\cdot\|$, i.e., $\|\mathbf{X}\|^\dagger = \sup_{\|\mathbf{X}'\| \leq 1} \langle \mathbf{X}, \mathbf{X}' \rangle$ for all $\mathbf{X} \in \mathbb{R}^{m \times n}$.

2.2 LINEAR MINIMIZATION ORACLE AND TRUST REGION

Let us look at the problem from the perspective of linear minimization oracle (lmo) and unconstrained stochastic conditional gradient descent (uSCG) (Pethick et al. (2025b)). lmo is defined as:

$$\text{lmo}(\mathcal{S}) \in \arg \min_{\mathbf{X} \in \mathcal{S}} \langle \mathbf{S}, \mathbf{X} \rangle, \quad (2)$$

where \mathcal{S} is some set. We are interested in the case when \mathcal{S} is a ball in our $\|\cdot\|$ norm:

$$\mathcal{S} := \mathcal{B}_\eta := \{\mathbf{X} \mid \|\mathbf{X}\| \leq \eta\}. \quad (3)$$

uSCG update is defined as: $\mathbf{X}^{k+1} = \mathbf{X}^k + \gamma_k \text{lmo}(\mathbf{M}^k)$, where $\mathbf{M}^{k+1} = (1 - \alpha_{k+1})\mathbf{M}^k + \alpha_{k+1}g(\mathbf{X}^k, \xi_k)$ is a momentum.

It can be easily shown that the formula is equivalent to

$$\mathbf{X}^{k+1} = \mathbf{X}^k - \gamma_k \eta \arg \max_{\mathbf{X} \in \mathcal{B}_1} \langle \mathbf{S}, \mathbf{X} \rangle = \mathbf{X}^k - \gamma_k \eta \{\Delta \in \mathcal{B}_1 \mid \langle \mathbf{M}^k, \Delta \rangle = \|\mathbf{M}^k\|^\dagger\}. \quad (4)$$

Let us set $\gamma_k = 1$. This transforms algorithm defined by eq. (4) into Algorithm 1 from Kovalev (2025). Therefore, we can view eq. (4) both as an lmo-based algorithm and as a trust-region algorithm.

3 DIFFERENT NORMS $\|\cdot\|$ IMPLY DIFFERENT UPDATES

Based on different norms $\|\cdot\|$, we simplify the update defined by the aforementioned equation:

$$\mathbf{X}^{k+1} = \mathbf{X}^k - \eta \{\Delta \in \mathcal{B}_1 \mid \langle \mathbf{M}^k, \Delta \rangle = \|\mathbf{M}^k\|^\dagger\} \quad (5)$$

In all the work, we define $\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}^\top$ as components of the singular value decomposition of \mathbf{M}^k : $\mathbf{M}^k = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$. We use common notations: $\mathbf{U} = [u_1, u_2, \dots, u_r]$, $\mathbf{\Sigma} = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r)$, and $\mathbf{V} = [v_1, v_2, \dots, v_r]$.

$\|\mathbf{M}^k\|_F$ and NSGD

Lemma 1. When $\|\cdot\| = \|\cdot\|_F$, eq. (5) turns into:

$$\mathbf{X}^{k+1} = \mathbf{X}^k - \eta \frac{\mathbf{M}^k}{\|\mathbf{M}^k\|_F} \quad (6)$$

It is an interesting observation, because in other works (Pethick et al. (2025b)), $\|\cdot\|_F$ was used to recover SGD. The difference is in how one states the problem.

$\|\mathbf{M}^k\|_{\text{op}}$ and **Muon**

Lemma 2. When $\|\cdot\| = \|\cdot\|_{\text{op}}$, eq. (5) turns into:

$$\mathbf{X}^{k+1} = \mathbf{X}^k - \eta UV^\top \quad (7)$$

$\|\mathbf{M}^k\|_{\text{nuc}}$ and **Neon**

Lemma 3. When $\|\cdot\| = \|\cdot\|_{\text{nuc}}$, eq. (5) turns into:

$$\mathbf{X}^{k+1} = \mathbf{X}^k - \eta u_1 v_1^\top \quad (8)$$

We name the derived algorithm *Neon*. In the section Matrix side of updates, we will discuss how to compute an update efficiently.

$\|\mathbf{M}^k\|_{F*}^\dagger$ and **F-Muon** We define $\|\cdot\|_{F*}$ as a convex combination of $\|\cdot\|_{\text{nuc}}$ and $\|\cdot\|_F$:

$$\|\mathbf{X}\|_{F*} = \alpha \|\mathbf{X}\|_{\text{nuc}} + (1 - \alpha) \|\mathbf{X}\|_F, \quad (9)$$

where $\alpha \in [0, 1]$ defines a specific norm of $F*$ -family.

Lemma 4. When $\|\cdot\| = \|\cdot\|_{F*}^\dagger$, eq. (5) turns into:

$$\mathbf{X}^{k+1} = \mathbf{X}^k - \eta (\alpha UV^\top + (1 - \alpha) \frac{\mathbf{M}^k}{\|\mathbf{M}^k\|_F}) \quad (10)$$

We name the derived algorithm *F-Muon*. It turns out that F-Muon is a convex combination of Normalized SGD and Muon, which is curious. The implications are significant and discussed in the following sections.

$\|\mathbf{M}^k\|_{F2}^\dagger$ and **F-Neon** We define $\|\cdot\|_{F2}$ as a convex combination of $\|\cdot\|_{\text{op}}$ and $\|\cdot\|_F$:

$$\|\mathbf{X}\|_{F2} = \alpha \|\mathbf{X}\|_{\text{op}} + (1 - \alpha) \|\mathbf{X}\|_F, \quad (11)$$

where $\alpha \in [0, 1]$ defines a specific norm of F2-family.

Lemma 5. When $\|\cdot\| = \|\cdot\|_{F2}^\dagger$, eq. (5) turns into:

$$\mathbf{X}^{k+1} = \mathbf{X}^k - \eta (\alpha u_1 v_1^\top + (1 - \alpha) \frac{\mathbf{M}^k}{\|\mathbf{M}^k\|_F}) \quad (12)$$

We name the derived algorithm *F-Neon*. It turns out that F-Neon is a convex combination of Normalized SGD and Neon, which is curious. The implications are significant and discussed in the following sections.

$\|\mathbf{M}^k\|_{\text{KF-k}}^\dagger$ and **Muon, Neon, and Dion without error feedback** We remind the reader that the Ky Fan k-norm Fan (1951), which we denote as $\|\cdot\|_{\text{KF-k}}$, is the sum of the k largest singular values of the matrix. It can be proven that $\|\cdot\|_{\text{KF-k}}^\dagger = \max\{\frac{1}{k} \|\cdot\|_{\text{nuc}}, \|\cdot\|_{\text{op}}\}$ (see Appendix). Special cases of Ky Fan k-norm are the Ky Fan 1-norm, which is the spectral norm, and the Ky Fan $\min\{m, n\}$ -norm, which is the nuclear norm.

Lemma 6. When $\|\cdot\| = \|M^k\|_{\text{KF-k}}^\dagger$, eq. (5) turns into:

$$\mathbf{X}^{k+1} = \mathbf{X}^k - \eta \sum_{i=1}^k u_i v_i^\top \quad (13)$$

This lemma recovers the updates of Muon, when $k = \min\{m, n\}$, of Neon, when $k = 1$, and of Dion Ahn & Xu (2025) without the error feedback.

Moreover, one can consider $\text{F-KF-k-norm} = \alpha \|\cdot\|_{\text{KF-k}} + (1-\alpha) \|\cdot\|_{\text{F}}$, the dual to which will produce algorithms like F-Dion without the error feedback.

3.1 ALGORITHMS FOR MATRICES \leftrightarrow ALGORITHMS FOR VECTORS

Table that compares Muons and Neons to vector algorithms (NSGD).

4 MATRIX SIDE OF UPDATES

4.1 THEORY

Not only formulas, but also citations! Solutions to efficiently find some parts of SVD:

- Basic SVD (complexity asymptotic!)
- Newton-Schulz (complexity asymptotic!)
- Our job: Power iterations, Randomized SVD and Lanczos

4.2 EXPERIMENTS

Nikolay's task: SVD, RSVD, Lanczos, Newton-Schulz. We use torch and cupy to test it. Options: all possible on torch, or additionally to compare all on cupy. The goal: to find how much do we lose due to cupy

5 TRUST REGION BOUNDS FOR L-SMOOTH FUNCTIONS

First, we analyze the problem in the unstochastic case. From Corollary 1 of Kovalev (2025), we directly get the following result that matches lower bounds, as was noted in that article.

Lemma 7. To reach the precision $\min_{k=1\dots K} \|\nabla f(\mathbf{X}_k)\|^\dagger \leq \varepsilon$ by the iterations equation 5 under the conditions of Assumption (A1), it is sufficient to choose the stepsize η and the number of iterations K as follows:

$$\eta = \mathcal{O}\left(\frac{\varepsilon}{L}\right), \quad K = \mathcal{O}\left(\frac{L\Delta_0}{\varepsilon^2}\right). \quad (14)$$

In the stochastic case, from Corollary of 2 of Kovalev (2025), we directly get the following result that pnce more matches lower bounds:

Lemma 8. To reach the precision $\mathbb{E} \min_{k=1\dots K} \|\nabla f(\mathbf{X}_k)\|^\dagger \leq \varepsilon$ by equation 5 under the assumptions Assumption (A2), Assumption (A1), Assumption (A3), it is sufficient to choose the parameters as follows:

$$\eta = \mathcal{O}\left(\min\left\{\frac{\varepsilon}{L}, \frac{\varepsilon^3}{\rho^2 \sigma^2 L}\right\}\right), \quad \alpha = \mathcal{O}\left(\min\left\{1, \frac{\varepsilon^2}{\rho^2 \sigma^2}\right\}\right), \quad (15)$$

$$K = \mathcal{O}\left(\max\left\{\frac{\rho \sigma}{\varepsilon}, \frac{\rho^3 \sigma^3}{\varepsilon^3}, \frac{L\Delta_0}{\varepsilon^2}, \frac{L\Delta_0 \rho^2 \sigma^2}{\varepsilon^4}\right\}\right). \quad (16)$$

As the norms $\|\cdot\|_{\text{F}}$, $\|\cdot\|_{\text{nuc}}$, $\|\cdot\|_{\text{F}}$ are almost proportional to each other when $m, n \rightarrow \infty$ (with high probability for random matrices), the expected convergence guarantees in terms of $\|\cdot\|_{\text{F}}$ are the same

(it can be easily shown by noting that $\|\mathbf{X}\| \sim \alpha \|\mathbf{X}\|_F$, $\|\nabla f(\mathbf{X})\|^\dagger \sim \frac{1}{\alpha} \|\nabla f(\mathbf{X})\|_F$, and expressing L -constant via L_F -constant for the Frobenius norm).

From the theory of random martices and the Marchenko-Pastur law, we get that random $\mathbf{M} \in \mathbb{R}^{m \times n}$: $\mathbf{M} \sim \mathcal{N}(0, \sigma^2 \mathbb{R}^{m \times n})$ has the following asymptotics of its norms:

Nuclear: $\sigma n \sqrt{m}$

Frobenius: $\sigma \sqrt{mn}$

Spectral: $\sigma(\sqrt{m} + \sqrt{n})$

This means that for square random matrices $n \times n$ the following asymptotics take place: $\|\cdot\|_F \sim \frac{\sqrt{n}}{2} \|\cdot\|_{\text{op}}$ and $\|\cdot\|_{\text{nuc}} \sim \frac{n}{2} \|\cdot\|_{\text{op}}$.

6 EXPERIMENTS

6.1 RANDOMIZED LINEAR LEAST SQUARES

Since the provided by other authors Kovalev (2025); Riabinin et al. (2025) theoretical guarantees are almost norm-independent, we have to test them in practice.

To test the bounds from Kovalev (2025) in practice, we construct the following L-smooth problem:

$$F(\mathbf{X}) = \frac{1}{2} \langle (\mathbf{X} - \mathbf{S}), \mathbf{M}(\mathbf{X} - \mathbf{S})\mathbf{N} \rangle \quad (17)$$

where $\mathbf{X} \in \mathbb{R}^{m \times n}$, $m = 10$, $n = 10$, $\mathbf{S} \in \mathbb{R}^{m \times n}$, $\mathbf{M} \in \mathbb{S}_+^m$ and $\mathbf{N} \in \mathbb{S}_+^n$ are positive-semidefinite matrices. Spectra of \mathbf{M} and \mathbf{N} are uniformly distributed in $(0, 1)$ interval.

It is easy to derive the gradient

$$\nabla F(\mathbf{X}) = \mathbf{M}(\mathbf{X} - \mathbf{S})\mathbf{N}, \quad (18)$$

Let us define γ as $\|\cdot\| \sim \gamma \|\cdot\|_F$, which is asymptotics from the previous section. Then $\|\cdot\|^\dagger \sim \frac{1}{\gamma} \|\cdot\|_F$, as $\|\cdot\|_F^\dagger = \|\cdot\|_F$. Hence, $\|\cdot\| \sim \gamma^2 \|\cdot\|^\dagger$

Then $\|\nabla f(\mathbf{X}) - \nabla f(\mathbf{Y})\|^\dagger = \|\mathbf{M}(\mathbf{X} - \mathbf{Y})\mathbf{N}\|^\dagger \leq \|\mathbf{M}\|^\dagger \|\mathbf{N}\|^\dagger \|\mathbf{X} - \mathbf{Y}\|^\dagger \sim \|\mathbf{M}\|^\dagger \|\mathbf{N}\|^\dagger \gamma^2 \|\mathbf{X} - \mathbf{Y}\|$, and $L \sim \gamma^2 \|\mathbf{M}\|^\dagger \|\mathbf{N}\|^\dagger$.

Now for the known norms:

Frobenius (NSGD): $L = \|\mathbf{M}\|_F \|\mathbf{N}\|_F$

Spectral (Muon): $\gamma \sim \frac{2}{\sqrt{n}} \implies L \sim \frac{2}{n} \|\mathbf{M}\|_{\text{nuc}} \|\mathbf{N}\|_{\text{nuc}}$

Nuclear (Neon): $\gamma \sim \sqrt{n} \implies L \sim n \|\mathbf{M}\|_{\text{op}} \|\mathbf{N}\|_{\text{op}}$

We take learning rate η and iteration number K from eq. (14).

We are interested to get $\varepsilon_2 = 1e-1$ precision in $\|\cdot\|_{\text{op}}$. Hence, we target $\varepsilon = \varepsilon_2$ for the nuclear norm (Neon), $\varepsilon = \frac{n}{2}$ for the spectral norm (Muon), and $\varepsilon = \frac{\sqrt{n}}{2}$ for the Frobenius norm (NSGD).

An independent experiment To make the problem more real-world, we do not use known information about smoothness. We run NormalizedSGD, Muon, F-Muon with $\alpha = 1/2$, Neon, and F-Neon with $\alpha = 1/2$ for 100 000 iterations with learning rate = 1e-3.

6.2 LOGISTIC REGRESSION

6.3 BENCHMARKS

CNN benchmark NanoGPT benchmark

7 RELATED WORK

As Muon Jordan et al. (2024) is a very successful and popular optimizer for functions of weight matrices, a lot of research has been put into, first, further improving its performance, and, second, in explaining its success.

Improvements of Muon. Regarding the first point, in less than a year, a large number of applications and improvements of Muon has been proposed. Liu et al. (2025) adapted the algorithm for training language models larger than NanoGPT. Shah et al. (2025) organized efficient hyperparameter transfer by combining Muon with maximal update parametrization. To construct their COSMOS optimizer, Chen et al. (2025) applied computationally intensive updates of SOAP optimizer to a low-dimensional “leading eigensubspace” while using memory-efficient methods like Muon for the remaining parameters. Amsel et al. (2025) proposed a more efficient alternative to Newton-Schulz operations. Si et al. (2025) introduced AdaMuon which combines element-wise adaptivity with orthogonal updates. We suppose that the described above or similar techniques can be applied to our optimizers as well. For example, F-Muon also benefits from faster alternatives to Newton-Schulz iterations, and Neon may be a great substitute to Muon in COSMOS, because, as we have shown in *the Matrix side of the updates*, Lanczos algorithms is much faster than Newton-Schulz iterations on large matrices.

Theory behind Muon. Regarding the second point, there has been a prolonged gap in theory behind Muon, simplistic derivation of Bernstein (2025) based on Bernstein & Newhouse (2024) excluded. This gap, as it seems to us, is not even now completely closed. For example, Kovalev (2025) has provided convergence guarantees of Muon in various settings, from which, however, Muon’s supremacy cannot be recovered. Indeed, although the obtained bounds depend on the norm choice, the asymptotics of the convergence remain the same as for NSGD and other optimizers, $K = \mathcal{O}(\varepsilon^{-4})$ in a L -smooth stochastic case.

Similar drawback has a recent article Riabinin et al. (2025), where L -smoothness assumption is replaced with a more practical (L_0, L_1) -smoothness. The authors derived from their theorems optimal stepsizes for Muon and Scion that match fine-tuned stepsizes reported by Pethick et al. (2025b). But still, they did not explain why, for instance, NSGD is inferior to Muon in training large-language models.

We suppose that the reason for the recorded by us discrepancy between Neon and Muon performance, both of which are described by Scion or Gluon frameworks, lies in the structure of the norm ball, which must be an object of further research.

The nuclear norm in lmo. As we found out only when writing this article, the nuclear norm has been already explored in the context of the linear minimization oracle. Pethick et al. (2025a) applied it to create ν SAM, a new sharpness-aware minimization technique. The update from their lemma 3.1 coincides with the update of Neon, but is used for completely different purposes. In addition, Pethick et al. (2025a) use power iterations to find u_1 and v_1^\top , while we suggest utilizing much more efficient and precise Lanczos algorithm.

The Ky Fan Norm and Dion. Rank- k Centralized Dion, Algorithm 1 from Ahn & Xu (2025), without an error feedback and scaling of the update, turns out to be an lmo-based algorithm under the $\|\cdot\|_{KF-k}^\dagger$ norm, which we described in *Different norms imply different updates*. Reported by the authors necessity of using error feedback to obtain satisfactory accuracy may take place in the cases of our algorithms as well, for example, F-Neon. This we leave to future research.

8 CONCLUSION

In this article, we have generalized several successful algorithms, like Muon and Dion, to the lmo-based algorithms in the $\|\cdot\|_{KF-k}^\dagger$ norm. Also we have proposed the technique of “regularizing” the updates with NSGD, a trick to increase the robustness of the algorithms and motivated by the consideration of the $\|\cdot\|_{F*}$ and $\|\cdot\|_{F2}$ norms. Generalizations of well-known norms and subsequent combination of them may further improve performance of lmo-based algorithms. If a theory is developed that explains the discrepancy between performance of different algorithms based on the matrix norms, one will probably be able to apply it to the generalizations and combinations of the

norms as well, which leaves an open space to construct norms specific for given architectures and probably even their parameters.

REFERENCES

- Kwangjun Ahn and Byron Xu. Dion: A communication-efficient optimizer for large models, 2025. URL <https://arxiv.org/abs/2504.05295>.
- Noah Amsel, David Persson, Christopher Musco, and Robert M Gower. The polar express: Optimal matrix sign methods and their application to the muon algorithm. *arXiv preprint arXiv:2505.16932*, 2025.
- Jeremy Bernstein. Deriving muon, 2025. URL <https://jeremybernste.in/writing/deriving-muon>.
- Jeremy Bernstein and Laker Newhouse. Old optimizer, new norm: An anthology. *arXiv preprint arXiv:2409.20325*, 2024.
- Weizhu Chen, Chen Liang, Tuo Zhao, Zixuan Zhang, Hao Kang, Liming Liu, Zichong Li, and Zhenghao Xu. Cosmos: A hybrid adaptive optimizer for memory-efficient training of llms, 2025. URL <https://arxiv.org/abs/2502.17410>.
- Ashok Cutkosky, Harsh Mehta, and Francesco Orabona. Momentum-based variance reduction in nonconvex sgd. *Advances in Neural Information Processing Systems*, 2020.
- Ky Fan. Maximum properties and inequalities for the eigenvalues of completely continuous operators. *Proceedings of the National Academy of Sciences*, 37(11):760–766, 1951.
- Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 2013.
- Saeed Ghadimi, Guanghui Lan, and Hongchao Zhang. Accelerated gradient methods for nonconvex nonlinear and stochastic programming. *Mathematical Programming*, 2016.
- Eduard Gorbunov, Dmitry Kovalev, and Peter Richtárik. Linearly converging error compensated sgd. *Advances in Neural Information Processing Systems*, 2020.
- Samuel Horváth, Dmitry Kovalev, and Peter Richtárik. Stochastic recursive momentum for nonconvex optimization. *arXiv preprint arXiv:2302.07731*, 2023.
- 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/>.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Dmitry Kovalev. Understanding gradient orthogonalization for deep learning via non-euclidean trust-region optimization. *arXiv preprint arXiv:2503.12645*, 2025.
- 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*, 2025.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Thomas Pethick, Parameswaran Raman, Lenon Minorics, Mingyi Hong, Shoham Sabach, and Volkan Cevher. ν SAM: Memory-efficient sharpness-aware minimization via nuclear norm constraints. *Transactions on Machine Learning Research*, 2025a. ISSN 2835-8856. URL <https://openreview.net/forum?id=V6ia5hWIMD>.

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*, 2025b.

Artem Riabinin, Egor Shulgin, Kaja Gruntkowska, and Peter Richtárik. Gluon: Making muon & scion great again!(bridging theory and practice of lmo-based optimizers for llms). *arXiv preprint arXiv:2505.13416*, 2025.

Ishaan Shah, Anthony M Polloreno, Karl Stratos, Philip Monk, Adarsh Chaluvvaraju, Andrew Hojel, Andrew Ma, Anil Thomas, Ashish Tanwer, Darsh J Shah, et al. Practical efficiency of muon for pretraining. *arXiv preprint arXiv:2505.02222*, 2025.

Chongjie Si, Debing Zhang, and Wei Shen. Adamuon: Adaptive muon optimizer. *arXiv preprint arXiv:2507.11005*, 2025.

Ruoyu Sun et al. Momentum methods for stochastic optimization: A survey and new results. *arXiv preprint arXiv:2302.06675*, 2023.

9 APPENDIX

9.1 FORMAL ASSUMPTIONS

Stochastic gradient estimator. We assume access to a stochastic estimator $g(\cdot; \xi): \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m \times n}$ of the gradient $\nabla f(\cdot)$, where $\xi \sim \mathcal{D}$ is a random variable sampled from a probability distribution \mathcal{D} . We assume that the stochastic gradient estimator $g(\cdot; \xi)$ is unbiased and has bounded variance, that is, the following relations hold:

$$\mathbb{E}_{\xi \sim \mathcal{D}} [g(\mathbf{X}; \xi)] = \nabla f(\mathbf{X}) \quad \text{and} \quad \mathbb{E}_{\xi \sim \mathcal{D}} [\|g(\mathbf{X}; \xi) - \nabla f(\mathbf{X})\|_F^2] \leq \sigma^2 \quad \text{for all } \mathbf{X} \in \mathbb{R}^{m \times n}, \quad (\text{A2})$$

where $\sigma > 0$ is a positive variance parameter, and $\|\cdot\|_F$ is the standard Euclidean, i.e. Frobenius, norm induced by the inner product $\langle \cdot, \cdot \rangle$, i.e., $\|\mathbf{X}\|_F = \sqrt{\langle \mathbf{X}, \mathbf{X} \rangle} = \sqrt{\text{tr}(\mathbf{X}^\top \mathbf{X})}$. These assumptions have been widely adopted for the analysis of many stochastic gradient optimization algorithms (Ghadimi & Lan, 2013; Ghadimi et al., 2016; Cutkosky et al., 2020; Sun et al., 2023; Horváth et al., 2023; Gorbunov et al., 2020).

It is important to highlight that while Assumption (A1) uses the dual norm $\|\cdot\|^\dagger$ to measure the difference between the gradients, the variance in Assumption (A2) is measured with respect to the Frobenius norm $\|\cdot\|_F^2$, which is necessary to properly utilize the unbiasedness property of the stochastic gradient estimator $g(\cdot; \xi)$. Therefore, we need to provide a connection between these norms using the following inequality:

$$\|\mathbf{X}\|^\dagger \leq \rho \cdot \|\mathbf{X}\|_F \quad \text{for all } \mathbf{X} \in \mathbb{R}^{m \times n}, \quad (\text{A3})$$

where $\rho > 0$ is a positive constant. Note that such a constant always exists due to the norm equivalence theorem, which always holds in the finite-dimensional space $\mathbb{R}^{m \times n}$. We recount ρ for different norms $\|\cdot\|$ in the appendix.

9.2 NORMS $\|\cdot\|_{F*}^\dagger$, $\|\cdot\|_{F2}^\dagger$, AND $\|\cdot\|_{KF-k}^\dagger$

Here we provide the derivation of these norms and plot them in the case of 2×2 matrices.

9.3 UPDATES DERIVATIONS

Proof of lemma 2 follows from eq. (10) with $\alpha = 1$. Indeed, $\|\cdot\|_{\text{op}}^\dagger = 1 \cdot \|\cdot\|_{\text{nuc}} + 0 \cdot \|\cdot\|_F$.

Proof of lemma 4: Since $\|\cdot\|^\dagger = \|\cdot\|_{F*}^\dagger = \|\cdot\|_{F*}$, the goal is to reach $\|\mathbf{M}^k\|_{F*} = \alpha \text{tr } \Sigma + (1 - \alpha) \|\mathbf{M}^k\|_F$.

Let us note that $\Delta = \alpha \mathbf{U} \mathbf{V}^\top + (1 - \alpha) \frac{\mathbf{M}^k}{\|\mathbf{M}^k\|_F}$ delivers this value. Indeed, by the trace property, $\langle \mathbf{M}^k, \Delta \rangle = \langle \mathbf{U} \Sigma \mathbf{V}^\top, \alpha \mathbf{U} \mathbf{V}^\top + (1 - \alpha) \frac{\mathbf{U} \Sigma \mathbf{V}^\top}{\|\mathbf{M}^k\|_F} \rangle = \alpha \text{tr } \Sigma + (1 - \alpha) \|\mathbf{M}^k\|_F = \|\mathbf{M}^k\|_{F*}$, which completes the proof.

Proof of lemma 3 follows from eq. (12) with $\alpha = 1$. Indeed, $\|\cdot\|_{\text{nuc}}^\dagger = 1 \cdot \|\cdot\|_{\text{op}} + 0 \cdot \|\cdot\|_{\text{F}}$.

Proof of lemma 5: Since $\|\cdot\|^\dagger = \|\cdot\|_{\text{F}^2}^{\dagger\dagger} = \|\cdot\|_{\text{F}^2}$, the goal is to reach $\|\mathbf{M}^k\|_{\text{F}^2} = \alpha\sigma_1 + (1 - \alpha)\|\mathbf{M}^k\|_{\text{F}}$.

Let us note that $\Delta = \alpha(u_1 v_1^\top) + (1 - \alpha) \frac{\mathbf{M}^k}{\|\mathbf{M}^k\|_{\text{F}}}$ delivers this value. Indeed, by the trace property and singular vectors orthogonality, $\langle \mathbf{M}^k, \Delta \rangle = \langle \mathbf{U} \Sigma \mathbf{V}^\top, \alpha u_1 v_1^\top + (1 - \alpha) \frac{\mathbf{U} \Sigma \mathbf{V}^\top}{\|\mathbf{M}^k\|_{\text{F}}} \rangle = \alpha \text{tr} \text{diag}(\sigma_1, 0, \dots, 0) + (1 - \alpha) \|\mathbf{M}^k\|_{\text{F}} = \|\mathbf{M}^k\|_{\text{F}^2}$, which completes the proof.

Proof of lemma 6: Since $\|\cdot\|^\dagger = \|\cdot\|_{\text{KF}-k}^{\dagger\dagger} = \|\cdot\|_{\text{KF}-k}$, the goal to reach $\|\mathbf{M}^k\|_{\text{KF}-k}$.

Let us note that $\Delta = \sum_{i=1}^k u_i v_i^\top$ delivers the value. Indeed, $\langle \mathbf{M}^k, \Delta \rangle = \langle \mathbf{U} \Sigma \mathbf{V}^\top, \sum_{i=1}^k u_i v_i^\top \rangle = \sum_{i,j=1}^{r,k} \langle u_i \sigma_i v_i^\top, u_j v_j^\top \rangle = \sum_{i=1}^k \sigma_i = \|\mathbf{M}^k\|_{\text{KF}-k}$, which completes the proof.

9.4 TECHNICAL DETAILS OF THE EXPERIMENTS

Note: nuSAM and their update. They had no Lanczos