

Delving into Muon and Beyond: Deep Analysis and Extensions

Xianbiao Qi^{1*} Marco Chen^{2*} Jiaquan Ye¹ Yelin He¹ Rong Xiao¹

Abstract

The Muon optimizer has recently attracted considerable attention for its strong empirical performance and use of orthogonalized updates on matrix-shaped parameters, yet its underlying mechanisms and relationship to adaptive optimizers such as Adam remain insufficiently understood. In this work, we aim to address these questions through a unified spectral perspective. Specifically, we view Muon as the $p = 0$ endpoint of a family of spectral transformations of the form $U\Sigma^pV^\top$, and consider additional variants with $p = \frac{1}{2}$, $p = \frac{1}{4}$, and $p = 1$. These transformations are applied to both first-moment updates, as in momentum SGD, and to root-mean-square (RMS) normalized gradient updates as in Adam. To enable efficient computation, we develop a coupled Newton iteration that avoids explicit singular value decomposition. Across controlled experiments, we find that RMS-normalized updates yield more stable optimization than first-moment updates. Moreover, while spectral compression provides strong stabilization benefits under first-moment updates, the Muon update ($p = 0$) does not consistently outperform Adam. These results suggest that Muon is best understood as an effective form of spectral normalization, but not a universally superior optimization method. Our source code will be released at <https://github.com/Ocram7/BeyondMuon>.

1. Introduction

Optimization methods (Robbins & Monro, 1951; Nesterov, 1983; Duchi et al., 2011; Hinton et al., 2012; Adam et al., 2014; Loshchilov & Hutter, 2019; Gupta et al., 2018; Jordan et al., 2024) are a cornerstone of modern deep learning. Among them, Adam (Adam et al., 2014) and its variants (Loshchilov & Hutter, 2019; Shazeer & Stern, 2018) are widely adopted due to their strong empirical stability

and robustness across architectures and scales. By combining momentum with variance-based normalization, Adam incorporates second-moment information while remaining computationally efficient, and has consequently become the default optimizer in large-scale neural network training (Touvron et al., 2023a;b; Dubey et al., 2024; Guo et al., 2025; Yang et al., 2025).

Recently, Muon (Jordan et al., 2024) has attracted growing attention as an alternative optimization strategy, particularly in large language models (Team et al., 2025). Unlike most conventional optimizers that operate element-wise, Muon applies matrix-level transformations to gradients, enforcing orthogonality through a spectral operation. Empirical evidence suggests that Muon can improve convergence behavior, leading to growing interest in the community. Despite growing adoption, the effectiveness and role of Muon remain insufficiently understood. Existing work has either focused on empirical performance gains (Jordan et al., 2024; Team et al., 2025) or on mathematical analysis (Newhouse, 2025; Bernstein, 2025a; Buchanan, 2025), but these studies are often accompanied by additional techniques such as QK-Norm or QK-Clip (Team et al., 2025). Consequently, it remains unclear to what extent the reported improvements can be attributed to Muon itself, how Muon relates to established adaptive optimizers such as Adam, and whether Muon can be systematically improved. These ambiguities motivate a careful and controlled re-examination of Muon.

In this work, our goal is to study the effectiveness of the Muon optimizer through a unified spectral framework and controlled empirical evaluation. We view Muon as the $p = 0$ endpoint of a family of spectral transformations of the form $U\Sigma^pV^\top$, and consider fractional variants with $p = \frac{1}{2}$ and $p = \frac{1}{4}$, as well as $p = 1$, which recovers a standard gradient update. We evaluate these transformations under both first-moment updates, as in momentum SGD, and root-mean-square-normalized updates, as in Adam. To make fractional spectral updates practical, we introduce a coupled Newton-Schulz iteration method that computes $U\Sigma^{\frac{1}{2}}V^\top$ and $U\Sigma^{\frac{1}{4}}V^\top$ using only matrix multiplications, avoiding explicit SVD.

To ensure fair and interpretable comparisons, we design our experiments to be highly controlled. We explicitly decouple matrix and vector learning rates for all methods, disable

^{*}Equal contribution ¹Intellifusion Inc. ²Tsinghua University. Correspondence to: Xianbiao Qi <qixianbiao@gmail.com>.

weight decay, and avoid auxiliary techniques such as QK-Norm and QK-Clip.

Under these settings, our experiments yield three key findings: (i) Muon ($p = 0$) is crucial for stabilizing first-moment updates and substantially improves robustness over mSGD; (ii) when applied to second-moment-normalized updates, Muon-like orthogonalization does not outperform Adam and appears inferior to partial spectral compression; and (iii) across spectral variants, Adam-style second-moment methods are consistently stronger than first-moment ones.

Our main contributions are summarized as follows:

- We introduce a unified spectral framework $U\Sigma^pV^\top$ that places Muon as the $p = 0$ endpoint and propose intermediate variants with $p = \frac{1}{2}$ and $p = \frac{1}{4}$.
- We develop a coupled Newton-Schulz iteration method that enables efficient computation of the fractional spectral updates $p = \frac{1}{2}$ and $p = \frac{1}{4}$ without explicit SVD.
- We provide a rigorous and controlled empirical comparison between Muon, Adam, and their spectral variants across first and second-moment-normalized updates, isolating spectral effects from other confounding techniques.

2. Related Work

Optimizers before Transformer. Before the Transformer (Vaswani et al., 2017), CNNs (LeCun et al., 2002; He et al., 2016) and shallow recurrent models such as LSTMs (Hochreiter & Schmidhuber, 1997) dominated neural network design. These architectures were typically of moderate depth with limited cumulative nonlinearity, resulting in relatively small variation in layer-wise Lipschitz constants (Qi et al., 2023a;b). Consequently, training with a global learning rate was generally effective. Under this regime, classical stochastic optimization methods (Bottou et al., 2018) were sufficient. Stochastic Gradient Descent (SGD) (Robbins & Monro, 1951) and its variants (Nesterov, 1983; Johnson & Zhang, 2013), often combined with simple learning rate schedules, were the standard choices.

Although adaptive methods such as Adagrad (Duchi et al., 2011), RMSProp (Hinton et al., 2012), and Adam (Adam et al., 2014) had been proposed, they were not widely adopted at the time. This can be attributed to two main factors. First, weight decay was primarily viewed as a regularization technique, and its interaction with adaptive optimizers was poorly understood; in particular, coupling weight decay with gradients rather than weights significantly degraded the performance of Adam-type methods (Loshchilov & Hutter, 2019). Second, in CNN-based models, inter-layer differences in Lipschitz constants were insufficient to necessitate per-parameter adaptivity, and momentum SGD

already achieved strong empirical performance. As a result, Adam did not consistently outperform mSGD in this setting.

Optimizers after Transformer. The introduction of Transformers (Vaswani et al., 2017) fundamentally altered the optimization landscape of deep learning. Compared to earlier architectures, Transformer-based models are substantially deeper and exhibit stronger nonlinearities (Shazeer, 2020), arising from attention mechanisms (Vaswani et al., 2017), residual connections (He et al., 2016), and normalization layers (Ba et al., 2016). Consequently, different layers and parameter groups often exhibit markedly different curvature and gradient statistics, often leading to training instability and suboptimal convergence.

AdamW (Loshchilov & Hutter, 2019) resolved the improper coupling between weight decay and adaptive learning rates in Adam, leading to significantly improved performance on Transformer models. As a result, AdamW consistently outperforms momentum SGD in large-scale Transformer training and has become a standard optimizer choice. Subsequent work has proposed a range of alternative optimizers, including Adafactor (Shazeer & Stern, 2018), SignSGD (Bernstein et al., 2018), LAMB (You et al., 2019), Adan (Xie et al., 2024), Lion (Chen et al., 2023), Sophia (Liu et al., 2024), Mars (Yuan et al., 2025), and related variants (Zhang et al., 2025; Liang et al., 2025). These methods primarily target improved stability, efficiency, or convergence in large-scale, highly nonlinear training regimes.

Matrix-based Optimizers. Matrix-based optimizers exploit the structured geometry of parameters by operating directly on matrix-valued variables, rather than treating parameters as independent scalars. Early work in this direction focused on making second-order optimization tractable via structured approximations. K-FAC (Martens & Grosse, 2015) introduced a Kronecker-factored approximation of the Fisher information matrix, enabling efficient layer-wise preconditioning, while Shampoo (Gupta et al., 2018) extended this idea by applying Kronecker factorizations along multiple matrix dimensions, leading to improved stability. These methods share the goal of incorporating curvature information through structured matrix approximations.

More recently, Muon (Jordan et al., 2024) has emerged as a matrix-based optimizer that departs from explicit curvature modeling. Muon directly updates matrix-valued parameters via orthogonalization and spectral transformations computed using Newton-Schulz iterations (Schulz, 1933; Higham, 2008), relying solely on matrix multiplications. Since its introduction, Muon has inspired a growing body of follow-up work (Team et al., 2025; Si et al., 2025; Frans et al., 2025; Su, 2025; Lau et al., 2025), including analyses that clarify its numerical linear algebra foundations (Bernstein, 2025a; Newhouse, 2025) and geometric interpretations that view it as implicit manifold-aware opti-

mization (Bernstein, 2025b; Buchanan, 2025).

Remarks. Unlike several previous works on optimizer benchmarking (Wen et al., 2025; Frans et al., 2025), our study is characterized by the following features: (1) we introduce a unified spectral-based formulation of the form $U\Sigma^pV^\top$ with $p \in [0, 1]$, and instantiate it with four concrete choices, $p \in \{1, \frac{1}{2}, \frac{1}{4}, 0\}$. For $U\Sigma^{\frac{1}{2}}V^\top$ and $U\Sigma^{\frac{1}{4}}V^\top$, we develop a coupled Newton–Schulz algorithm to solve them. (2) we experimentally show that matrix-based optimizers with second-moment information are stronger than their pure first-moment counterparts. (3) we design our experiments to be highly controlled. We explicitly decouple matrix and vector learning rates for all methods, disable weight decay, and avoid auxiliary techniques. This allows us to isolate the intrinsic capabilities of each optimizer, enabling a fair and transparent comparison.

3. Delving into Muon and Beyond

3.1. Baselines: Adam and Muon

As in Muon, our study focuses exclusively on *matrix-shaped parameters* $\mathbf{W} \in \mathbb{R}^{m \times n}$; all vector parameters are optimized using standard Adam. In this subsection, we briefly review Adam and Muon and fix notation.

Adam. Adam (Adam et al., 2014) is a popular stochastic optimizer that combines momentum with adaptive, second-moment-based normalization. Given the gradient matrix $\mathbf{G}_t = \nabla_{\mathbf{W}_t} \ell(\mathbf{W}_t)$ at iteration t , Adam maintains exponentially decaying estimates of the first and second moments:

$$\mathbf{M}_t = \beta_1 \mathbf{M}_{t-1} + (1 - \beta_1) \mathbf{G}_t, \quad (1)$$

$$\mathbf{V}_t = \beta_2 \mathbf{V}_{t-1} + (1 - \beta_2) \mathbf{G}_t^2, \quad (2)$$

and updates parameters via

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta_t \mathbf{M}_t \oslash \sqrt{\mathbf{V}_t}. \quad (3)$$

Here, β_1 and β_2 are the first- and second-moment decay rates, η_t is the learning rate, and \oslash denotes elementwise division. From a preconditioning perspective, Adam applies a diagonal matrix $\mathbf{V}_t^{-1/2}$ to the momentum \mathbf{M}_t , normalizing each coordinate by an estimate of its second-moment and yielding bounded, scale-adaptive updates. For brevity, throughout the rest of this work we refer to $\mathbf{M}_t \oslash \sqrt{\mathbf{V}_t}$ as the *RMS-normalized* update, since \mathbf{V}_t tracks an exponential moving average of squared gradients.

Muon. Muon (Jordan et al., 2024) replaces Adam’s elementwise normalization with a *spectral normalization* of the update matrix. Given the momentum matrix \mathbf{M}_t with singular value decomposition $\mathbf{M}_t = \mathbf{U}\Sigma\mathbf{V}^\top$, Muon defines the update direction as the polar factor

$$\mathcal{P}(\mathbf{M}_t) := \mathbf{U}\mathbf{V}^\top,$$

which discards the singular values and retains only the left and right singular vectors. Computing the polar factor is typically referred to as matrix orthogonalization.

From an optimization perspective, orthogonalization enforces *unit magnitude along every singular direction*. In contrast to Adam, which normalizes updates on an elementwise basis via a diagonal preconditioner, Muon equalizes the strength of updates across all directions of the gradient matrix. As a result, Muon removes anisotropic scaling in the spectrum of \mathbf{M}_t , yielding directionally balanced but magnitude-agnostic updates.

Newton–Schulz approximation. Computing the polar factor $\mathcal{P}(\mathbf{M}_t) = \mathbf{U}\mathbf{V}^\top = \mathbf{M}_t(\mathbf{M}_t^\top \mathbf{M}_t)^{-1/2}$ via an explicit SVD is prohibitively expensive for large models. However, we can approximate the inverse square root using *Newton–Schulz iteration*, which relies only on matrix multiplications. Let $\mathbf{A} = \alpha \mathbf{M}_t^\top \mathbf{M}_t$, where the scaling constant α is chosen such that $\|\mathbf{I} - \mathbf{A}\|_2 \leq 1$. The iteration

$$\mathbf{Z}_{k+1} = \frac{1}{2} \mathbf{Z}_k (3\mathbf{I} - \mathbf{A}\mathbf{Z}_k^2), \quad \mathbf{Z}_0 = \mathbf{I}, \quad (4)$$

converges quadratically to $\frac{1}{\sqrt{\alpha}}(\mathbf{M}_t^\top \mathbf{M}_t)^{-1/2}$. After K iterations, the polar factor is approximated as

$$\mathcal{P}(\mathbf{M}_t) \approx \sqrt{\alpha} \mathbf{M}_t \mathbf{Z}_K. \quad (5)$$

In practice, Muon (Jordan et al., 2024) employs Algorithm 2 in Appendix B, which is a slightly more direct and aggressive version of the Newton–Schulz iteration shown in Equation 4. The official implementation directly computes the polar factor with a procedure that converges for initializations where the singular values are adequately normalized.

3.2. Momentum vs. RMS-Normalized Updates

Where to apply spectral transforms. Muon (Jordan et al., 2024) applies its spectral operation to the *first-moment* momentum \mathbf{M}_t . However, in large-scale Transformer training, Adam-style optimizers that incorporate second-moment statistics generally outperform momentum SGD-based optimizers that only leverage first-moment information. This observation motivates us to examine whether spectral transformations yield additional benefits when applied to rms-normalized updates. Concretely, we study the two inputs

$$\mathbf{O}_t^{\text{mom}} := \mathbf{M}_t \quad \text{and} \quad \mathbf{O}_t^{\text{rms}} := \mathbf{M}_t \oslash \sqrt{\mathbf{V}_t},$$

which correspond to the first-moment update shown in Equation 1 and the rms-normalized update from Equation 3.

How $\mathbf{O}_t^{\text{mom}}$ and $\mathbf{O}_t^{\text{rms}}$ differ. The key distinction between these two regimes lies in how update magnitudes are controlled. In first-order momentum methods, the update \mathbf{M}_t aggregates gradients over time without explicit normalization by their scale. As a result, the norm of the update can

grow with accumulated gradient magnitude, especially under high variance or anisotropic curvature. Consequently, the update scale is not bounded by the optimizer.

In contrast, RMS-normalized updates explicitly rescale the momentum using second-moment statistics. The update $\mathbf{M}_t \oslash \sqrt{\mathbf{V}_t}$ normalizes each coordinate by an estimate of its raw second moment, yielding an update whose magnitude is provably bounded. For Adam-style methods, this bound, $\frac{1-\beta_1}{\sqrt{1-\beta_2}}$ (Qi et al., 2023b), depends only on the decay rates and is independent of the raw gradient magnitude. This intrinsic normalization leads to improved numerical stability and more predictable optimization dynamics.

3.3. A Spectral Family of Transformations

We introduce a family of spectral gradient transformations that enables a unified analysis of Muon-inspired spectral methods. Given SVD, $\mathbf{O}_t = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$, we define the general spectral transformation

$$\Psi_p(\mathbf{O}_t) := \mathbf{U}\mathbf{\Sigma}^p\mathbf{V}^\top, \quad p \in [0, 1]. \quad (6)$$

The exponent p directly controls how the singular spectrum is rescaled. Writing the SVD in its rank-one form, we obtain

$$\Psi_p(\mathbf{O}_t) = \sum_{i=1}^d \sigma_i^p \mathbf{u}_i \mathbf{v}_i^\top, \quad (7)$$

which makes it explicit that decreasing p progressively compresses the singular values $\{\sigma_i\}$ while preserving the singular directions $\{\mathbf{u}_i, \mathbf{v}_i\}$.

This formulation subsumes several existing methods and enables systematic comparisons across spectral behaviors. In particular, $p = 1$ leaves the input unchanged, recovering standard methods like mSGD and Adam, while $p = 0$ maps all nonzero singular values to 1, yielding an orthogonalized update direction. When the input is the first-moment momentum $\mathbf{O}_t = \mathbf{O}_t^{\text{mom}}$, $p = 0$ corresponds to Muon.

3.4. Instantiating the Spectral Family

From a spectral operator to concrete optimizers. We instantiate our spectral update family by applying the map $\Psi_p(\cdot)$ from Equation 6 to one of two matrix-valued inputs: the first-moment momentum $\mathbf{O}_t^{\text{mom}}$, or the RMS-normalized update $\mathbf{O}_t^{\text{rms}}$. We consider four spectral exponents

$$p \in \left\{1, \frac{1}{2}, \frac{1}{4}, 0\right\},$$

where $p = 1$ leaves the input unchanged and $p = 0$ collapses the spectrum to the polar factor. The intermediate powers $p = \frac{1}{2}$ and $p = \frac{1}{4}$ provide milder alternatives to Muon’s fully flattened spectrum, yielding a controlled interpolation

$$\mathbf{U}\mathbf{\Sigma}^1\mathbf{V}^\top \rightarrow \mathbf{U}\mathbf{\Sigma}^{1/2}\mathbf{V}^\top \rightarrow \mathbf{U}\mathbf{\Sigma}^{1/4}\mathbf{V}^\top \rightarrow \mathbf{U}\mathbf{\Sigma}^0\mathbf{V}^\top.$$

Combining two inputs with four exponents produces the eight optimizer instances evaluated in this work.

Naming convention. Each instance is identified by its *input family* and *spectral exponent*. We write

$$\text{BASEX} : \Delta \mathbf{W}_t \propto \Psi_p(\mathbf{O}_t),$$

where $\text{BASE} \in \{\mathbf{mSGD}, \mathbf{Adam}\}$ selects $\mathbf{O}_t^{\text{mom}}$ or $\mathbf{O}_t^{\text{rms}}$, and the suffix $X \in \{\emptyset, \mathbf{S}, \mathbf{Q}, \mathbf{Z}\}$ encodes the exponent

$$\begin{aligned} \emptyset &\leftrightarrow p = 1 \text{ (identity)}, & \mathbf{S} &\leftrightarrow p = \frac{1}{2} \text{ (square-root)}, \\ \mathbf{Q} &\leftrightarrow p = \frac{1}{4} \text{ (quarter-power)}, & \mathbf{Z} &\leftrightarrow p = 0 \text{ (zero-power / polar)}. \end{aligned}$$

Momentum-input family ($\mathbf{O}_t^{\text{mom}} = \mathbf{M}_t$). Under the momentum input, the four instances are

$$\mathbf{mSGD}, \mathbf{mSGDS}, \mathbf{mSGDQ}, \mathbf{mSGDZ},$$

corresponding to $p \in \{1, \frac{1}{2}, \frac{1}{4}, 0\}$ respectively. The endpoint **mSGDZ** is equivalent to **Muon** since $\Psi_0(\mathbf{M}_t) = \mathcal{P}(\mathbf{M}_t) = \mathbf{U}\mathbf{V}^\top$.

RMS-normalized-input family ($\mathbf{O}_t^{\text{rms}} = \mathbf{M}_t \oslash \sqrt{\mathbf{V}_t}$). Under the RMS-normalized input, the four instances are

$$\mathbf{Adam}, \mathbf{AdamS}, \mathbf{AdamQ}, \mathbf{AdamZ},$$

again corresponding to $p \in \{1, \frac{1}{2}, \frac{1}{4}, 0\}$ respectively.

Spectral anisotropy control. The four spectral variants in each input family share the same singular vectors (\mathbf{U}, \mathbf{V}) and differ only in how the singular values are rescaled. Thus, the effect of the spectral transformation can be understood purely as reshaping the anisotropy of the singular spectrum. To quantify anisotropy, we use the spectral condition number

$$\kappa(\mathbf{O}_t) = \frac{\sigma_{\max}}{\sigma_{\min}},$$

which measures the spread of the singular values. When $\kappa(\mathbf{O}_t)$ is large, the input is highly ill-conditioned and a small number of dominant singular directions can disproportionately influence the update.

If the input has singular values $\{\sigma_i\}_{i=1}^d$, then Ψ_p produces singular values $\{\sigma_i^p\}_{i=1}^d$, so the condition number becomes

$$\kappa_p = \frac{\sigma_{\max}^p}{\sigma_{\min}^p} = \kappa(\mathbf{O}_t)^p.$$

Therefore, decreasing p reduces spectral anisotropy:

$$\kappa_1 \geq \kappa_{1/2} \geq \kappa_{1/4} \geq \kappa_0 = 1,$$

Indeed, for the toy spectrum in Figure 1, $\kappa_1 = 90000$, $\kappa_{1/2} = 300$, $\kappa_{1/4} = \sqrt{300}$, and $\kappa_0 = 1$. This viewpoint clarifies how the exponent p continuously interpolates between the original spectrum ($p = 1$) and complete spectral flattening ($p = 0$).

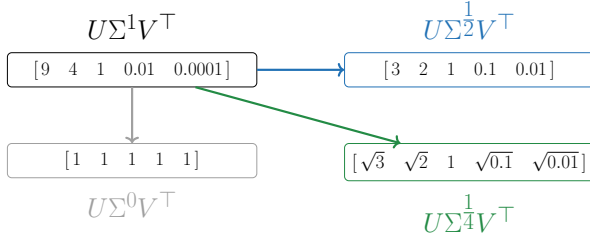


Figure 1. Effect of spectral exponent p on the singular spectrum (illustrative numbers). Decreasing p compresses the spectrum: large singular values are damped relative to small ones, and $p = 0$ maps all nonzero singular values to 1.

The effect of spectral compression. As shown in Figure 1, decreasing p progressively compresses the *entire* singular-value spectrum. Hence, directions associated with large singular values are attenuated, while those associated with small singular values (less than one) are amplified.

3.5. Efficient Computation

Computing $\Psi_{1/2}(\mathbf{O}_t)$ and $\Psi_{1/4}(\mathbf{O}_t)$. A practical challenge in spectral methods is computational efficiency, as direct SVD computation is infeasible for large models. Fortunately, we can rewrite $\Psi_{1/2}(\mathbf{O}_t)$ in a form that avoids explicit SVD. Let $\mathbf{O}_t = \mathbf{U}\Sigma\mathbf{V}^\top$ and assume $m \geq n$. Then

$$\Psi_{1/2}(\mathbf{O}_t) = \mathbf{U}\Sigma^{1/2}\mathbf{V}^\top = \mathbf{O}_t(\mathbf{O}_t^\top\mathbf{O}_t)^{-1/4}.$$

Thus, it suffices to compute an inverse fourth root of the symmetric matrix $\mathbf{X} := \mathbf{O}_t^\top\mathbf{O}_t \in \mathbb{R}^{n \times n}$. We can compute $\mathbf{X}^{1/2}$ and $\mathbf{X}^{-1/2}$ efficiently using the coupled Newton-Schulz algorithm discussed below, and obtain $\mathbf{X}^{-1/4}$ by applying the same procedure to $\mathbf{X}^{1/2}$ (i.e., computing $(\mathbf{X}^{1/2})^{-1/2} = \mathbf{X}^{-1/4}$). Consequently, $\Psi_{1/2}(\mathbf{O})$ can be implemented using only matrix multiplications as $\Psi_{1/2}(\mathbf{O}) = \mathbf{O}\mathbf{X}^{-1/4}$. A similar matrix-multiplication-only procedure can be used for the quarter-power spectral transformations.

Coupled Newton-Schulz Iteration. Coupled Newton-Schulz iteration (Higham, 2008) provides an efficient and numerically stable procedure for simultaneously computing the matrix square root $\mathbf{X}^{1/2}$ and its inverse $\mathbf{X}^{-1/2}$.

Starting from the initialization $\mathbf{Y}_0 = \mathbf{X}$ and $\mathbf{Z}_0 = \mathbf{I}$, the method applies a coupled update that repeatedly refines both quantities using only matrix multiplications. Specifically, each iteration is defined as

$$\begin{aligned} \mathbf{Y}_{k+1} &= \frac{1}{2} \mathbf{Y}_k (3\mathbf{I} - \mathbf{Z}_k \mathbf{Y}_k), \\ \mathbf{Z}_{k+1} &= \frac{1}{2} (3\mathbf{I} - \mathbf{Z}_k \mathbf{Y}_k) \mathbf{Z}_k, \end{aligned} \quad (8)$$

which symmetrically updates \mathbf{Y}_k and \mathbf{Z}_k through the shared correction term $3\mathbf{I} - \mathbf{Z}_k \mathbf{Y}_k$. When appropriately

Algorithm 1 Coupled Newton-Schulz for $\mathbf{X}^{1/2}$ and $\mathbf{X}^{-1/2}$

Require: \mathbf{X} , number of iterations K

Ensure: $\mathbf{X}^{1/2}$, $\mathbf{X}^{-1/2}$

```

1:  $\mathbf{Y}_0 \leftarrow \mathbf{X}$ ,  $\mathbf{Z}_0 \leftarrow \mathbf{I}$ 
2:  $\alpha \leftarrow \|\mathbf{X}\|_F$ 
3:  $\mathbf{Y}_0 \leftarrow \mathbf{Y}_0/\alpha$ ,  $k \leftarrow 0$ 
4: while  $k < K$  do
5:    $\mathbf{T}_k \leftarrow 3\mathbf{I} - \mathbf{Z}_k \mathbf{Y}_k$ 
6:    $\mathbf{Y}_{k+1} \leftarrow \frac{1}{2} \mathbf{Y}_k \mathbf{T}_k$ 
7:    $\mathbf{Z}_{k+1} \leftarrow \frac{1}{2} \mathbf{T}_k \mathbf{Z}_k$ 
8:    $k \leftarrow k + 1$ 
9: end while
10: return  $\sqrt{\alpha} \mathbf{Y}_K$ ,  $\frac{1}{\sqrt{\alpha}} \mathbf{Z}_K$ 
```

scaled, this coupled iteration converges quadratically, driving $\mathbf{Y}_k \rightarrow \mathbf{X}^{1/2}$ and $\mathbf{Z}_k \rightarrow \mathbf{X}^{-1/2}$ simultaneously.

In practice, a normalization step based on the Frobenius norm of \mathbf{X} is applied at initialization to ensure numerical stability, and the final iterates are rescaled accordingly. In contrast to the standard one-sequence Newton-Schulz iteration, which can be obtained by eliminating one variable from the coupled Newton iteration under a commuting initialization, the coupled formulation simultaneously evaluates \mathbf{Y}_k and \mathbf{Z}_k , making it a more suitable tool for computing matrix roots and inverse roots. Algorithm 1 depicts our method.

4. Experiments

4.1. Experimental Settings

Setup. We conduct our experiments on nanoGPT¹ (Karpathy, 2022), a lightweight GPT-2 training codebase. We follow the standard GPT-2 configuration: GELU activations and a byte-pair encoding tokenizer with vocabulary size 50,257. Our GPT-2 model has 124M parameters. All runs use sequence length 1024, global batch size 480, and are trained on OpenWebText for 200K optimization steps with warmup for 2,000 steps. We do not use QK-Norm or QK-Clip. Across runs, we keep all settings fixed and vary only the optimizer and its learning rate.

For Muon, we use the reference implementation² from Jordan et al. (2024), which applies the matrix update $\mathbf{W}_{t+1} = \mathbf{W}_t - \eta_t \sqrt{\frac{\text{fan-out}}{\text{fan-in}}} \mathbf{U}\mathbf{V}^\top$, where $\mathbf{U}\mathbf{V}^\top$ is the polar factor computed via Algorithm 2.

Controlled comparison. To ensure fair and interpretable comparisons, we isolate the optimizer’s effect on *matrix-shaped* parameters and, crucially, *decouple matrix and vec-*

¹<https://github.com/karpathy/nanoGPT>

²<https://github.com/KellerJordan/Muon>

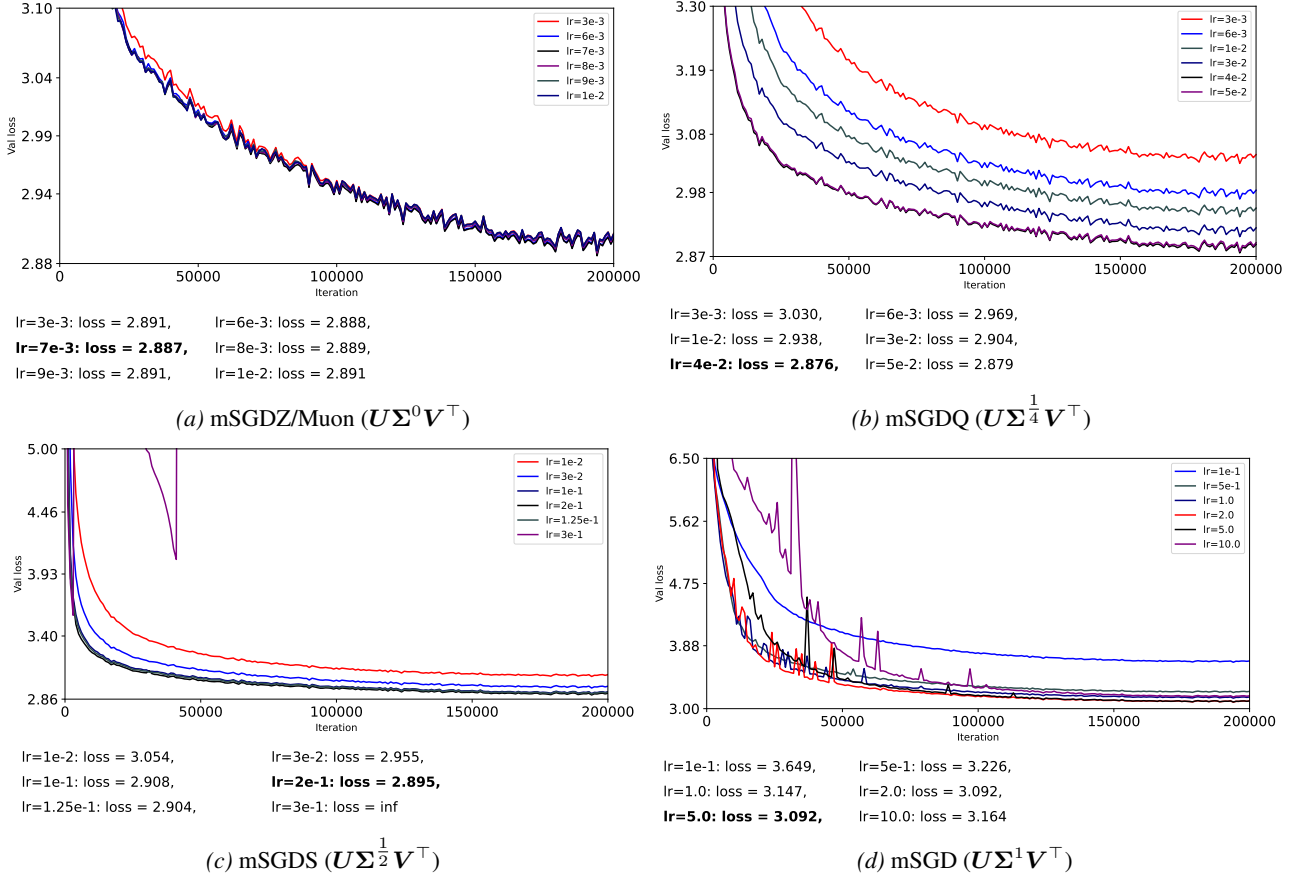


Figure 2. Overall comparison across four optimizers (mSGDZ/Muon, mSGDQ, mSGDS and mSGD) based on the **first-moment momentum** M_t . Each subfigure corresponds to a different optimizer.

tor learning rates for *all* optimizers.

Muon is typically evaluated with separate learning rates for matrix and vector parameters, whereas standard Adam baselines often share a single global learning rate. To eliminate this mismatch, we follow Muon and apply Adam to all vector-valued parameters for every method using a fixed learning rate $\text{lr}_{\text{vec}} = 3 \times 10^{-4}$. For matrix parameters, each optimizer is tuned with its own learning rate lr_{mat} .

We also disable weight decay for all runs ($\text{wd} = 0$). In standard implementations, weight decay is coupled to the learning rate through the update $\mathbf{W}_{t+1} = \mathbf{W}_t - \eta_t \Delta \mathbf{W}_t - \eta_t \lambda \mathbf{W}_t$, so jointly tuning (η_t, λ) introduces an additional degree of freedom that complicates controlled cross-optimizer comparisons. In addition, we avoid auxiliary stabilization techniques such as QK-Norm and QK-Clip. All remaining hyperparameters are held fixed across methods; in particular, we set $\beta_1 = 0.9$ and $\beta_2 = 0.95$ throughout.

Overall, this design decouples the matrix learning rate from other hyperparameters, including vector learning rates and weight decay, so performance differences are more directly attributable to the matrix-parameter update rule itself.

4.2. Learning Rate Tuning Protocol

We tune the matrix-parameter learning rate lr_{mat} for each optimizer using a two-stage procedure: a coarse logarithmic sweep to identify a stable scale, followed by a local refinement within that scale.

Coarse search. We first sweep lr_{mat} on a logarithmic grid spanning several orders of magnitude to locate the region where the optimizer transitions from divergence to effective learning. For instance, for AdamW we evaluate $\{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$. This stage intends to identify the correct order of magnitude rather than the exact optimum. Since most candidates in this sweep diverge, we do not plot the resulting curves.

Fine-grained search. After identifying a promising scale, we refine the search by evaluating a denser set of learning rates above and below the best coarse candidate. We stop refining once the selected learning rate performs better than its immediate neighbors in this local grid, indicating a stable local optimum within the explored range. For each optimizer, this refinement evaluates roughly ten candidates; for

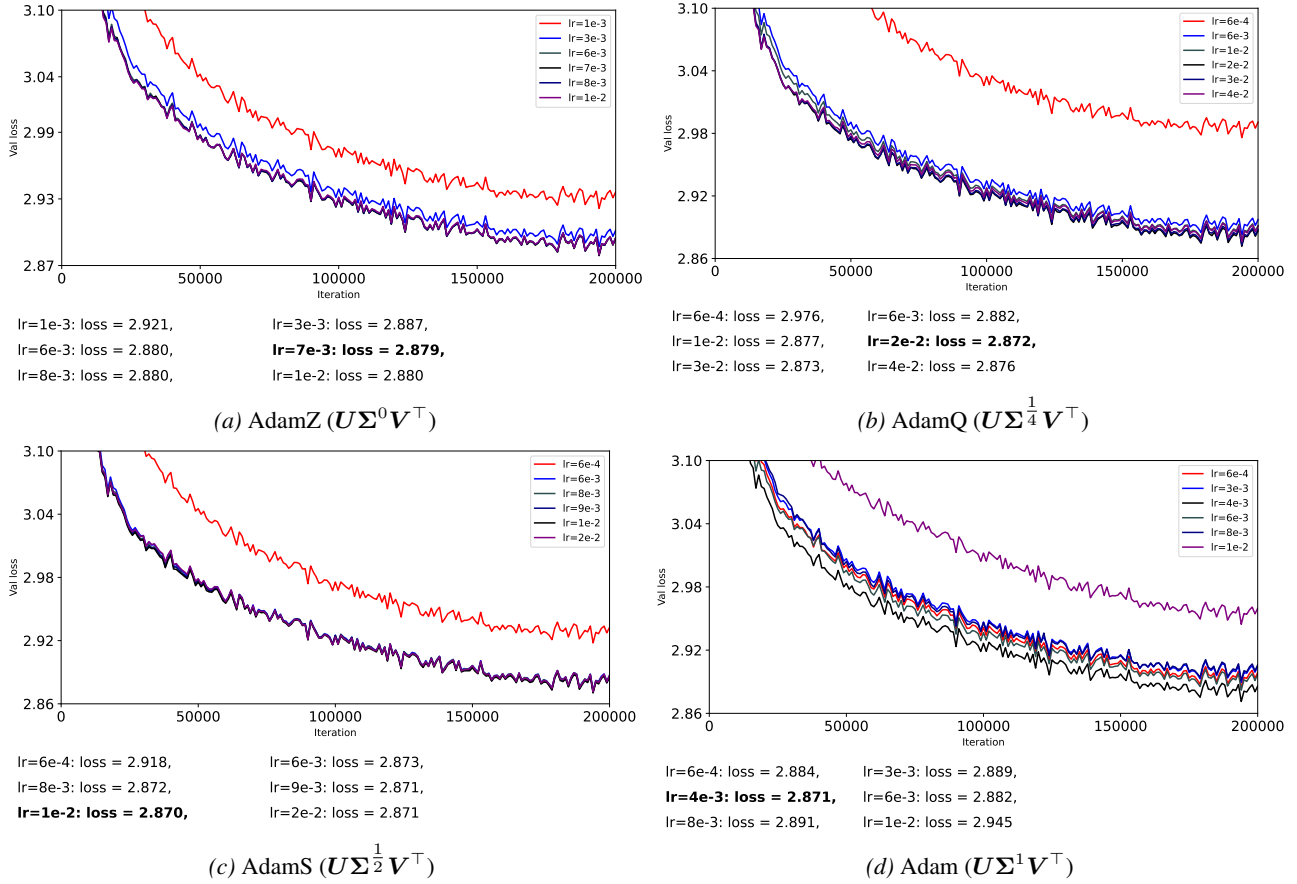


Figure 3. Overall comparison across four optimizers (AdamZ, AdamS, AdamQ and Adam) based on the **second-moment-normalized update** $M_t \oslash \sqrt{V_t}$. Each subfigure corresponds to a different optimizer.

readability, we plot six representative runs in the figures.

4.3. Effect of Spectral Exponent on Momentum Inputs

Figure 2 reports representative tuning results when the input O_t^{mom} is used. We make three observations:

- *Muon stabilizes momentum updates.* Muon (mSGDZ) is considerably more stable than mSGD across a wide learning-rate range. The best learning rate we find for Muon, 7×10^{-3} , closely matches the value reported in prior benchmarking (Wen et al., 2025).
- *Partial compression improves stability but not to Muon’s level.* The intermediate variants mSGDS and mSGDQ are more stable than mSGD, but remain less robust than Muon under aggressive learning rates.
- *After stability is achieved, moderate compression can outperform flattening.* Among the four momentum-input variants, mSGDQ attains the strongest performance once it trains stably, exceeding Muon and obtaining strong results at learning rates of 4×10^{-2} and 5×10^{-2} .

4.4. Effect of Spectral Exponent on Normalized Inputs

Similarly, Figure 3 shows our selected results given rms-normalized momentum O_t^{rms} as input. We observe that:

- *Spectral transforms yield smaller gains when the input is already normalized.* Applying spectral compression to O_t^{rms} produces limited improvements: AdamS can modestly broaden the stable learning-rate range, whereas stronger compression (AdamQ) offers little benefit and full flattening (AdamZ) degrades performance.
- *All four variants behave similarly at their best settings.* The peak performance differences among Adam, AdamS, and AdamQ are modest, suggesting that elementwise normalization already controls much of the harmful anisotropy that spectral compression targets.

4.5. Momentum vs. RMS-Normalized Inputs

In Figure 4, we compare the best-tuned run from each of the eight configurations. Two trends are clear. First, switching the input from the momentum O_t^{mom} to the RMS-normalized

update O_t^{rms} consistently improves both stability and final performance across spectral exponents. In particular, AdamZ (Muon-style orthogonalization applied to O_t^{rms}) is stronger than mSGDZ (the original Muon applied to O_t^{mom}). Second, within the RMS-normalized family, performance differences across p are relatively small, with AdamS emerging as a strong and stable choice.

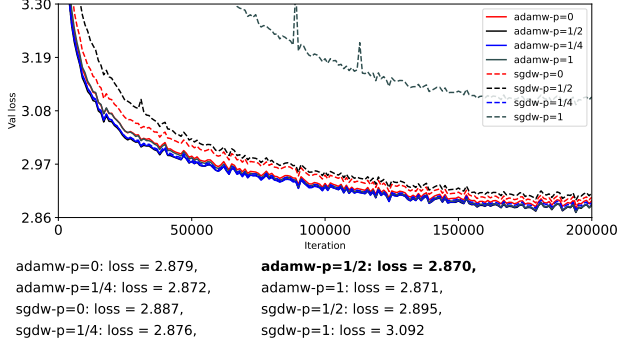


Figure 4. Comparison of four spectral exponents $p \in \{0, \frac{1}{4}, \frac{1}{2}, 1\}$ applied to either the first-moment momentum M_t or the RMS-normalized update $M_t \oslash \sqrt{V_t}$, each shown at its tuned learning rate. Dashed curves correspond to momentum-input variants, while solid curves correspond to RMS-normalized-input variants.

5. Discussion: Our Understanding of Muon

The primary goal of this work is to evaluate when Muon-style orthogonalization (and, more broadly, spectral compression) is a reliable optimization strategy. Accordingly, we interpret Muon’s empirical behavior in Section 4 through three honest conclusions.

I. Muon exhibits no significant performance edge over AdamW. In our controlled setting, Muon-style orthogonalization does not outperform Adam, whether applied to first-moment momentum or to RMS-normalized updates. This contrasts with prior reports of Muon converging faster than Adam (Jordan et al., 2024). We attribute the discrepancy primarily to the methodological differences noted in Subsection 4.1. Broadly, we employ universal learning-rate decoupling, avoid auxiliary techniques like QK-Norm and QK-Clip, and disable weight decay.

II. Muon is stable but aggressive. On first-moment momentum inputs, Muon provides a clear stabilization effect where orthogonalization significantly improves robustness over mSGD across most learning rates (Figure 2). This behavior aligns with Muon’s defining operation: given $G_t = U\Sigma V^\top$, Muon replaces the spectrum by Σ^0 , producing the polar factor UV^\top . By discarding singular-value magnitudes, Muon prevents directions with extremely large singular values from dominating the update, mitigating ill-conditioning through *direction-wise* spectral scale control.

However, this stabilization mechanism is also somewhat crude. Muon applies the same singular-value scaling to every direction, independent of curvature or signal strength. When some directions are legitimately larger than others, flattening them to unity can effectively reduce the step taken along informative directions. In this sense, Muon stabilizes by construction, but it does not selectively preserve meaningful anisotropy when such anisotropy is beneficial.

III. Muon may magnify noisy gradient directions.

Muon’s biggest limitation follows from the same flattening: by mapping all nonzero singular values to 1, it removes magnitude-based filtering. Directions corresponding to small singular values are no longer suppressed relative to dominant directions and can receive comparable update weight. This is especially problematic in the later stages of training where update matrices typically have low effective rank, and small-singular-value components may be primarily dominated by noise.

This effect is consistent with our observations that aggressive Muon-like orthogonalization on RMS-normalized inputs does not outperform Adam and is actually often worse than partial compression. One plausible explanation is that RMS normalization already bounds and regularizes per-coordinate scales, so additional full-spectrum flattening primarily reweights the update toward directions that RMS normalization alone would keep small.

Overall. Together, our results suggest that spectral compression can provide genuine stabilization benefits, especially when the input update is *not* already normalized. However, full orthogonalization ($p = 0$) is not always desirable: it can under-step along genuinely informative dominant modes (crudeness) and over-step along low-signal modes (noise magnification). Hence, our experiments do not support orthogonalization as a universally superior replacement for modern RMS-normalized second-moment optimizers.

6. Conclusion and Limitations

Conclusion. This paper studies matrix-based optimizers from a spectral perspective. We interpret a range of optimization methods under a unified spectral framework. Our results indicate that while spectral transformations like Muon have strong stabilizing properties, they cannot directly replace second-moment adaptive optimizers such as Adam.

Limitations and future work. Our study omits weight decay to enable a controlled comparison and relies on a coupled Newton–Schulz iteration procedure that can incur nontrivial overhead. Investigating how to combine matrix-based updates with regularization techniques such as weight decay and improving the practical efficiency of the coupled Newton–Schulz algorithm are directions for future work.

Impact Statement

“This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.”

References

- Adam, K. D. B. J. et al. A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 1412(6), 2014.
- Ba, J. L., Kiros, J. R., and Hinton, G. E. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- Bernstein, J. Deriving muon, 2025a. URL <https://jeremybernste.in/writing/deriving-muon>.
- Bernstein, J. Modular manifolds. *Thinking Machines Lab: Connectionism*, 2025b. doi: 10.64434/tml.20250926. <https://thinkingmachines.ai/blog/modular-manifolds/>.
- Bernstein, J., Wang, Y.-X., Azizadenesheli, K., and Anandkumar, A. signsgd: Compressed optimisation for non-convex problems. In *International Conference on Machine Learning*, pp. 560–569. PMLR, 2018.
- Bottou, L., Curtis, F. E., and Nocedal, J. Optimization methods for large-scale machine learning. *SIAM review*, 60(2):223–311, 2018.
- Buchanan, S. A faster manifold muon with ADMM. <https://sdbuchanan.com/blog/manifold-muon/>, 2025.
- Chen, X., Liang, C., Huang, D., Real, E., Wang, K., Liu, Y., Pham, H., Dong, X., Luong, T., Hsieh, C.-J., Lu, Y., and Le, Q. V. Symbolic discovery of optimization algorithms, 2023. URL <https://arxiv.org/abs/2302.06675>.
- Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al. The llama 3 herd of models. *arXiv e-prints*, pp. arXiv–2407, 2024.
- Duchi, J., Hazan, E., and Singer, Y. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- Frans, K., Abbeel, P., and Levine, S. What really matters in matrix-whitening optimizers? *arXiv preprint arXiv:2510.25000*, 2025.
- Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Gupta, V., Koren, T., and Singer, Y. Shampoo: Pre-conditioned stochastic tensor optimization. In *International Conference on Machine Learning*, pp. 1842–1850. PMLR, 2018.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Higham, N. J. Computing the polar decomposition—with applications. *SIAM Journal on Scientific and Statistical Computing*, 7(4):1160–1174, 1986.
- Higham, N. J. Computing real square roots of a real matrix. *Linear Algebra and its applications*, 88:405–430, 1987.
- Higham, N. J. *Functions of matrices: theory and computation*. SIAM, 2008.
- Hinton, G., Srivastava, N., and Swersky, K. Neural networks for machine learning lecture 6a overview of mini-batch gradient descent. *Cited on*, 14(8):2, 2012.
- Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- Horn, R. A. and Johnson, C. R. *Matrix analysis*. Cambridge university press, 2012.
- Johnson, R. and Zhang, T. Accelerating stochastic gradient descent using predictive variance reduction. *Advances in neural information processing systems*, 26, 2013.
- Jordan, K., Jin, Y., Boza, V., You, J., Cesista, F., Newhouse, L., and Bernstein, J. Muon: An optimizer for hidden layers in neural networks, 2024. URL <https://kellerjordan.github.io/posts/muon/>.
- Karpathy, A. NanoGPT. <https://github.com/karpathy/nanoGPT>, 2022.
- Lau, T. T.-K., Long, Q., and Su, W. Polargrad: A class of matrix-gradient optimizers from a unifying preconditioning perspective. *arXiv preprint arXiv:2505.21799*, 2025.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 2002.
- Liang, K., Chen, L., Liu, B., and Liu, Q. Cautious optimizers: Improving training with one line of code, 2025. URL <https://arxiv.org/abs/2411.16085>.
- Liu, H., Li, Z., Hall, D., Liang, P., and Ma, T. Sophia: A scalable stochastic second-order optimizer for language model pre-training, 2024. URL <https://arxiv.org/abs/2305.14342>.

- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization, 2019. URL <https://arxiv.org/abs/1711.05101>.
- Martens, J. and Grosse, R. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pp. 2408–2417. PMLR, 2015.
- Nesterov, Y. A method for unconstrained convex minimization problem with the rate of convergence $O(1/k^2)$. In *Doklady an ussr*, volume 269, pp. 543–547, 1983.
- Newhouse, L. Understanding muon, 2025. URL <https://lakernewhouse.com/muon>.
- Qi, X., Wang, J., Chen, Y., Shi, Y., and Zhang, L. Lipsformer: Introducing lipschitz continuity to vision transformers. *arXiv preprint arXiv:2304.09856*, 2023a.
- Qi, X., Wang, J., and Zhang, L. Understanding optimization of deep learning via jacobian matrix and lipschitz constant. *arXiv preprint arXiv:2306.09338*, 2023b.
- Robbins, H. and Monro, S. A stochastic approximation method. *The annals of mathematical statistics*, pp. 400–407, 1951.
- Schulz, G. Iterative berechnung der reziproken matrix. *ZAMM-Journal of Applied Mathematics and Mechanics/Zeitschrift für Angewandte Mathematik und Mechanik*, 13(1):57–59, 1933.
- Shazeer, N. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- Shazeer, N. and Stern, M. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, pp. 4596–4604. PMLR, 2018.
- Si, C., Zhang, D., and Shen, W. Adamuon: Adaptive muon optimizer, 2025. URL <https://arxiv.org/abs/2507.11005>.
- Su, W. Isotropic curvature model for understanding deep learning optimization: Is gradient orthogonalization optimal? *arXiv preprint arXiv:2511.00674*, 2025.
- Team, K., Bai, Y., Bao, Y., Chen, G., Chen, J., Chen, N., Chen, R., Chen, Y., Chen, Y., Chen, Y., et al. Kimi k2: Open agentic intelligence. *arXiv preprint arXiv:2507.20534*, 2025.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Wen, K., Hall, D., Ma, T., and Liang, P. Fantastic pretraining optimizers and where to find them. *arXiv preprint arXiv:2509.02046*, 2025.
- Xie, X., Zhou, P., Li, H., Lin, Z., and Yan, S. Adan: Adaptive nesterov momentum algorithm for faster optimizing deep models, 2024. URL <https://arxiv.org/abs/2208.06677>.
- Yang, A., Li, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Gao, C., Huang, C., Lv, C., Zheng, C., Liu, D., Zhou, F., Huang, F., Hu, F., Ge, H., Wei, H., Lin, H., Tang, J., Yang, J., Tu, J., Zhang, J., Yang, J., Yang, J., Zhou, J., Zhou, J., Lin, J., Dang, K., Bao, K., Yang, K., Yu, L., Deng, L., Li, M., Xue, M., Li, M., Zhang, P., Wang, P., Zhu, Q., Men, R., Gao, R., Liu, S., Luo, S., Li, T., Tang, T., Yin, W., Ren, X., Wang, X., Zhang, X., Ren, X., Fan, Y., Su, Y., Zhang, Y., Zhang, Y., Wan, Y., Liu, Y., Wang, Z., Cui, Z., Zhang, Z., Zhou, Z., and Qiu, Z. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- You, Y., Li, J., Reddi, S., Hseu, J., Kumar, S., Bhojanapalli, S., Song, X., Demmel, J., Keutzer, K., and Hsieh, C.-J. Large batch optimization for deep learning: Training bert in 76 minutes. *arXiv preprint arXiv:1904.00962*, 2019.
- Yuan, H., Liu, Y., Wu, S., Zhou, X., and Gu, Q. Mars: Unleashing the power of variance reduction for training large models, 2025. URL <https://arxiv.org/abs/2411.10438>.
- Zhang, Y., Chen, C., Li, Z., Ding, T., Wu, C., Kingma, D. P., Ye, Y., Luo, Z.-Q., and Sun, R. Adam-mini: Use fewer learning rates to gain more, 2025. URL <https://arxiv.org/abs/2406.16793>.
- Zur, J. Matrix functions. Technical report, Technische Universität Berlin, 2025. URL https://www.tu.berlin/fileadmin/www/40000110/Dauer-Wimis/Jan_Zur/matrix_functions_zur.pdf. Lecture note.

Almost all derivations are about functions of matrices, all these derivatinsos can be found in the following materials (Higham, 1986; 1987; 2008; Zur, 2025; Horn & Johnson, 2012).

A. Derivation of UV^\top based on Newton-Schulz method

Let the singular value decomposition of a matrix $M \in \mathbb{R}^{m \times n}$ be

$$M = USV^\top,$$

where $U \in \mathbb{R}^{m \times r}$, $V \in \mathbb{R}^{n \times r}$, $S = \text{Diag}(\sigma_1, \dots, \sigma_r)$, $\sigma_i > 0$, and U, V have orthonormal columns. Let us assume $m \leq n$

Our goal is to compute the orthogonal (polar) factor UV^\top *without explicitly performing an SVD*, using the Newton–Schulz iteration.

Step 1: Express UV^\top as a matrix function of M . We start from quantities that can be formed directly from M . Consider $M^\top M \in \mathbb{R}^{n \times n}$. Substituting the SVD of M gives

$$M^\top M = (USV^\top)^\top (USV^\top) = VSU^\top USV^\top = VS^2V^\top,$$

where we used $U^\top U = I_r$.

Hence, the inverse square root of $M^\top M$ (on its rank- r subspace) is

$$(M^\top M)^{-1/2} = VS^{-1}V^\top.$$

Left-multiplying by M yields

$$M(M^\top M)^{-1/2} = (USV^\top)(VS^{-1}V^\top) = U(SS^{-1})V^\top = UV^\top.$$

This identity shows that UV^\top can be obtained by removing the singular-value scaling of M via $(M^\top M)^{-1/2}$.

Step 2: Reduction to an inverse square root problem. Define

$$A := M^\top M \succeq 0.$$

The problem is now reduced to computing $A^{-1/2}$ using only matrix multiplication, without SVD.

Step 3: Newton–Schulz iteration for the inverse square root. Let us consider a scalar $a > 0$, consider the problem of computing $z = a^{-1/2}$. This is equivalent to solving $f(z) = \frac{1}{z^2} - a = 0$. Applying Newton’s method gives

$$z^+ = z - \frac{f(z)}{f'(z)} = z - \frac{z^{-2} - a}{-2z^{-3}} = \frac{1}{2}z(3 - az^2).$$

This is the scalar Newton–Schulz update.

Replacing the scalar a by a matrix $A \succeq 0$, the scalar z by a matrix Z , and scalar multiplication by matrix multiplication yields the matrix Newton–Schulz iteration

$$Z_{k+1} = \frac{1}{2}Z_k(3I - AZ_k^2).$$

This iteration is equal to Equation 4 in the main body. It involves only matrix multiplication, addition, and scalar scaling. After sufficient iterations,

$$Z_k \approx A^{-1/2}.$$

Step 4: Recovering UV^\top . Using the identity $UV^\top = MA^{-1/2}$, $A = M^\top M$, and the approximation $Z_k \approx A^{-1/2}$, we obtain

$$UV^\top \approx MZ_k.$$

B. Newton-Schulz method used by [Jordan et al. \(2024\)](#)

Algorithm 2 Newton-Schulz to solve UV^\top where $M = U\Sigma V^\top$

Require: M, K

Ensure: UV^\top

```

1: Coefficients:  $(a, b, c) \leftarrow (3.4445, -4.7750, 2.0315)$ 
2:  $Z_0 \leftarrow M$ 
3:  $\alpha \leftarrow \|Z_0\|_F$ 
4:  $Z_0 \leftarrow Z_0/\alpha$ ,  $k \leftarrow 0$ 
5: while  $k < K$  do
6:    $Z_{k+1} \leftarrow aZ_k + bZ_kZ_k^\top Z_k + c(Z_kZ_k^\top)^2 Z_k$ 
7:    $k \leftarrow k + 1$ 
8: end while
9: return  $Z_K$ 
    
```

C. Derive $X^{\frac{1}{2}}$ based on Coupled Newton-Schulz method

We aim to compute the matrix square root $X^{1/2}$. Let the unknown matrix Y satisfy

$$f(Y) := Y^2 - X = 0.$$

This defines a matrix equation $f : \mathbb{C}^{n \times n} \rightarrow \mathbb{C}^{n \times n}$.

Fréchet Derivative. For an arbitrary perturbation H , we have

$$f(Y + H) = (Y + H)^2 - X = Y^2 - X + (YH + HY) + H^2.$$

Therefore, the Fréchet derivative of f at Y is

$$Df(Y)[H] = YH + HY.$$

Definition of the Newton Step. The Newton increment ΔY at Y is defined as the solution of

$$Df(Y)[\Delta Y] = -f(Y),$$

that is, the Sylvester equation

$$Y\Delta Y + \Delta Y Y = -(Y^2 - X). \tag{9}$$

The update is then given by

$$Y^+ = Y + \Delta Y.$$

Up to this point, the derivation corresponds to the *exact matrix Newton method*, without any approximation.

A Coupled Newton System. The main difficulty of (9) is that it requires solving a Sylvester equation. A classical technique is to couple the unknown square root with its inverse, thereby converting the Sylvester solve into matrix multiplications.

Let Z approximate Y^{-1} (and ultimately $X^{-1/2}$). Consider the coupled system

$$\begin{cases} F_1(Y, Z) := Y^2 - X = 0, \\ F_2(Y, Z) := ZY - I = 0. \end{cases} \quad (10)$$

The exact solution is

$$(Y_*, Z_*) = (X^{1/2}, X^{-1/2}).$$

Fréchet Derivative of the Coupled System For perturbations $(\Delta Y, \Delta Z)$:

- For F_1 ,

$$DF_1(Y, Z)[\Delta Y, \Delta Z] = Y\Delta Y + \Delta Y Y.$$

- For F_2 ,

$$(Z + \Delta Z)(Y + \Delta Y) - I = (ZY - I) + (Z\Delta Y + \Delta Z Y) + (\Delta Z)(\Delta Y),$$

hence

$$DF_2(Y, Z)[\Delta Y, \Delta Z] = Z\Delta Y + \Delta Z Y.$$

Newton Linear System (Exact) The Newton step $(\Delta Y, \Delta Z)$ is defined by

$$DF(Y, Z)[\Delta Y, \Delta Z] = -F(Y, Z),$$

that is,

$$\begin{aligned} Y\Delta Y + \Delta Y Y &= -(Y^2 - X), \\ Z\Delta Y + \Delta Z Y &= -(ZY - I). \end{aligned} \quad (11)$$

This system is still an *exact Newton system*.

A Closed-Form Newton Step via a Symmetric Multiplicative Ansatz We adopt the symmetric multiplicative ansatz

$$Y^+ = YQ, \quad Z^+ = QZ, \quad (12)$$

where $Q = Q(Y, Z)$ is to be determined.

Define

$$R := ZY.$$

Then

$$R^+ := Z^+ Y^+ = (QZ)(YQ) = Q(ZY)Q = QRQ.$$

Hence, choosing Q reduces to the problem: given R , construct Q such that $QRQ \approx I$.

Deriving $Q = \frac{1}{2}(3I - ZY)$ from Newton's Principle. To approximate $R^{-1/2}$, ideally,

$$Q = R^{-1/2} \Rightarrow QRQ = I.$$

Consider the matrix function

$$\phi(R) := R^{-1/2}.$$

When R is close to the identity, write

$$R = I + E, \quad \|E\| \ll 1.$$

According to the Taylor expansion, we have

$$(I + E)^{-1/2} = I - \frac{1}{2}E + O(\|E\|^2).$$

Substituting $E = R - I$, we obtain

$$R^{-1/2} \approx I - \frac{1}{2}(R - I) = \frac{1}{2}(3I - R).$$

Thus define

$$Q := \frac{1}{2}(3I - R) = \frac{1}{2}(3I - ZY). \quad (13)$$

Deriving $Q = \frac{1}{2}(3I - ZY)$ from Newton's Principle Coupled Newton–Schulz Iteration.

Substituting (13) into (12), we obtain

$$\begin{aligned} Y^+ &= \frac{1}{2} Y(3I - ZY), \\ Z^+ &= \frac{1}{2} (3I - ZY) Z. \end{aligned}$$

Write it as an iteration equation, we have

$$\begin{aligned} Y_{k+1} &= \frac{1}{2} Y_k(3I - Z_k Y_k), & Y_0 &= X, \\ Z_{k+1} &= \frac{1}{2} (3I - Z_k Y_k) Z_k, & Z_0 &= I. \end{aligned} \quad (14)$$

This is the *coupled Newton–Schulz iteration*, which simultaneously drives

$$Y_k \rightarrow X^{1/2}, \quad Z_k \rightarrow X^{-1/2}.$$

D. Derive $X^{\frac{1}{4}}$ based on Coupled Newton-Schulz method

Algorithm 3 Two-times Coupled Newton-Schulz to solve $X^{\frac{1}{4}}$ and $X^{\frac{-1}{4}}$

Require: X, K

Ensure: $X^{\frac{1}{4}}, X^{\frac{-1}{4}}$

```

1:  $Y_0 \leftarrow X, Z_0 \leftarrow I$ 
2:  $\alpha \leftarrow \|X\|_F$ 
3:  $Y_0 \leftarrow \frac{Y_0}{\alpha}, k \leftarrow 0$ 
4: while  $k < K$  do
5:    $T_k \leftarrow 3I - Z_k Y_k$ 
6:    $Y_{k+1} \leftarrow \frac{1}{2} Y_k T_k$ 
7:    $Z_{k+1} \leftarrow \frac{1}{2} T_k Z_k$ 
8:    $k \leftarrow k + 1$ 
9: end while
10:  $Y_0 \leftarrow \sqrt{\alpha} Y_K, Z_0 \leftarrow I$ 
11:  $\beta \leftarrow \|Y_0\|_F$ 
12:  $Y_0 \leftarrow \frac{Y_0}{\beta}, k \leftarrow 0$ 
13: while  $k < K$  do
14:    $T_k \leftarrow 3I - Z_k Y_k$ 
15:    $Y_{k+1} \leftarrow \frac{1}{2} Y_k T_k$ 
16:    $Z_{k+1} \leftarrow \frac{1}{2} T_k Z_k$ 
17:    $k \leftarrow k + 1$ 
18: end while
19: return  $\sqrt{\beta} Y_K, \frac{1}{\sqrt{\beta}} Z_K$ 

```

E. Detailed results for different settings

Table 1. Hyperparameter sweeps for Muon.

Origin	$\Psi_p(O^{\text{mom}})$	Learning Rate	Weight Decay	Val Loss
First order	$U\Sigma^1 V^\top$	1e-1	0	3.649
First order	$U\Sigma^1 V^\top$	5e-1	0	3.226
First order	$U\Sigma^1 V^\top$	1.0	0	3.147
First order	$U\Sigma^1 V^\top$	2.0	0	3.092
First order	$U\Sigma^1 V^\top$	5.0	0	3.092
First order	$U\Sigma^1 V^\top$	10.0	0	3.164
First order	$U\Sigma^{\frac{1}{2}} V^\top$	1e-2	0	3.054
First order	$U\Sigma^{\frac{1}{2}} V^\top$	3e-2	0	2.955
First order	$U\Sigma^{\frac{1}{2}} V^\top$	1e-1	0	2.908
First order	$U\Sigma^{\frac{1}{2}} V^\top$	2e-1	0	2.895
First order	$U\Sigma^{\frac{1}{2}} V^\top$	1.25e-1	0	2.904
First order	$U\Sigma^{\frac{1}{2}} V^\top$	3e-1	0	inf
First order	$U\Sigma^{\frac{1}{4}} V^\top$	3e-3	0	3.030
First order	$U\Sigma^{\frac{1}{4}} V^\top$	6e-3	0	2.969
First order	$U\Sigma^{\frac{1}{4}} V^\top$	1e-2	0	2.938
First order	$U\Sigma^{\frac{1}{4}} V^\top$	3e-2	0	2.904
First order	$U\Sigma^{\frac{1}{4}} V^\top$	4e-2	0	2.876
First order	$U\Sigma^{\frac{1}{4}} V^\top$	5e-2	0	2.879
First order	$U\Sigma^0 V^\top$	3e-3	0	2.891
First order	$U\Sigma^0 V^\top$	6e-3	0	2.888
First order	$U\Sigma^0 V^\top$	7e-3	0	2.887
First order	$U\Sigma^0 V^\top$	8e-3	0	2.889
First order	$U\Sigma^0 V^\top$	9e-3	0	2.891
First order	$U\Sigma^0 V^\top$	1e-2	0	2.891

Table 2. Hyperparameter sweeps for Muon.

Origin	$\Psi_p(O^{\text{rms}})$	Learning Rate	Weight Decay	Val Loss
Second order	$U\Sigma^1 V^\top$	6e-4	0	2.884
Second order	$U\Sigma^1 V^\top$	3e-3	0	2.889
Second order	$U\Sigma^1 V^\top$	4e-3	0	2.871
Second order	$U\Sigma^1 V^\top$	6e-3	0	2.882
Second order	$U\Sigma^1 V^\top$	8e-3	0	2.891
Second order	$U\Sigma^1 V^\top$	1e-2	0	2.945
Second order	$U\Sigma^{\frac{1}{2}} V^\top$	6e-4	0	2.918
Second order	$U\Sigma^{\frac{1}{2}} V^\top$	6e-3	0	2.873
Second order	$U\Sigma^{\frac{1}{2}} V^\top$	8e-3	0	2.872
Second order	$U\Sigma^{\frac{1}{2}} V^\top$	9e-3	0	2.871
Second order	$U\Sigma^{\frac{1}{2}} V^\top$	1e-2	0	2.870
Second order	$U\Sigma^{\frac{1}{2}} V^\top$	2e-2	0	2.871
Second order	$U\Sigma^{\frac{1}{4}} V^\top$	6e-4	0	2.976
Second order	$U\Sigma^{\frac{1}{4}} V^\top$	6e-3	0	2.882
Second order	$U\Sigma^{\frac{1}{4}} V^\top$	1e-2	0	2.877
Second order	$U\Sigma^{\frac{1}{4}} V^\top$	2e-2	0	2.872
Second order	$U\Sigma^{\frac{1}{4}} V^\top$	3e-2	0	2.873
Second order	$U\Sigma^{\frac{1}{4}} V^\top$	4e-2	0	2.876
Second order	$U\Sigma^0 V^\top$	1e-3	0	2.921
Second order	$U\Sigma^0 V^\top$	3e-3	0	2.887
Second order	$U\Sigma^0 V^\top$	6e-3	0	2.880
Second order	$U\Sigma^0 V^\top$	7e-3	0	2.879
Second order	$U\Sigma^0 V^\top$	8e-3	0	2.880
Second order	$U\Sigma^0 V^\top$	1e-2	0	2.880