

ROOT: Robust Orthogonalized Optimizer for Neural Network Training

Wei He^{*1} Kai Han^{*1} Hang Zhou¹ Hanting Chen¹ Zhicheng Liu¹ Xinghao Chen¹ Yunhe Wang¹

¹Huawei Noah's Ark Lab

{hewei142, kai.han, zhouhang25, chenhanting, liuzhicheng15, xinghao.chen, yunhe.wang}@huawei.com

Abstract

The optimization of large language models (LLMs) remains a critical challenge, particularly as model scaling exacerbates sensitivity to algorithmic imprecision and training instability. The recent advances in optimizer improve convergence efficiency through momentum orthogonalization but suffers from two key robustness limitations: dimensional fragility in orthogonalization precision and vulnerability to outlier-induced noise. To address these robustness challenges, we introduce ROOT, a **Robust Orthogonalized OpTimizer** that enhances training stability through dual robustness mechanisms. First, we develop a dimension-robust orthogonalization scheme using adaptive Newton iterations with fine-grained coefficients tailored to specific matrix sizes, ensuring consistent precision across diverse architectural configurations. Second, we introduce an optimization-robust framework via proximal optimization that suppresses outlier noise while preserving meaningful gradient directions. Extensive experiments demonstrate that ROOT achieves significantly improved robustness, with faster convergence and superior final performance compared to both Muon and Adam-based optimizers, particularly in noisy and non-convex scenarios. Our work establishes a new paradigm for developing robust and precise optimizers capable of handling the complexities of modern large-scale model training. The code will be available at <https://github.com/huawei-noah/noah-research/tree/master/ROOT>.

1. Introduction

The escalating computational demands of pre-training Large Language Models (LLMs) (Brown et al., 2020; OpenAI,

2023; Chen et al., 2025) have positioned the design of optimization algorithms as a critical research frontier. The choice of an optimizer profoundly influences not only the convergence speed and final performance but also the substantial economic cost of model development. While Stochastic Gradient Descent (SGD) (Robbins & Monro, 1951) established the foundational paradigm, the need for greater efficiency and stability in navigating complex, high-dimensional loss landscapes has driven the development of adaptive methods. The overarching goal is to design optimizers that are computationally efficient, robust, and scalable, enabling faster training of increasingly powerful models.

Adam (Kingma, 2014) and its variant AdamW (Loshchilov & Hutter, 2017) have long been the de facto standards for training deep learning models. By incorporating momentum and per-parameter adaptive learning rates, Adam often achieves faster convergence than SGD. However, as model sizes surge into the billions of parameters, inherent limitations of adaptive methods become increasingly apparent, particularly numerical instability in mixed-precision training. Recent optimizers, such as Muon (Jordan et al., 2024; Liu et al., 2025), represent architectural shifts aimed at addressing these challenges. For instance, Muon views weight matrices as holistic entities rather than independent scalars, employing a Newton-Schulz iteration to approximate orthogonal factors from momentum matrices. Such approaches have demonstrated benefits in training speed and memory efficiency compared to AdamW.

Despite these advances, our analysis identifies two critical limitations prevalent among those commonly-used optimizers. First, methods relying on fixed-coefficient iterative approximations (*e.g.*, Newton-Schulz (Jordan et al., 2024)) often exhibit a substantial precision gap, particularly for certain matrix dimensions where approximation errors are markedly higher. This stems from a one-size-fits-all approach to numerical stabilization, which fails to adapt to the diverse spectral properties of weight matrices across different layers. Second, many adaptive optimizers show heightened sensitivity to outlier-induced gradient noise, which is a common phenomenon in large-scale training where anomalous samples produce gradient components with dis-

^{*}Equal contribution ¹Huawei Noah's Ark Lab.

proportionately large magnitudes. This can corrupt update directions and destabilize training.

To address these limitations, we introduce **ROOT (Robust Orthogonalized OpTimizer)**, a method designed to systematically enhance robustness from two key perspectives. First, to achieve *algorithmic robustness* against structural uncertainties, we replace the fixed-coefficient Newton-Schulz iteration with an adaptive scheme employing fine-grained, dimension-specific coefficients. This ensures high-precision orthogonalization across all layers, making the process robust to variations in matrix dimensions. Second, for *optimization robustness* against data-level noise, we incorporate a proximal optimization term that suppresses outlier-induced gradient noise via soft-thresholding, thereby stabilizing training without compromising convergence. Together, these contributions form a unified approach that significantly improves the robustness and reliability of the optimization process.

The main contributions of this work are summarized as follows:

- We identify two key robustness limitations in the orthogonalization-based optimizers: lack of algorithmic robustness due to imprecise orthogonalization across varying matrix dimensions, and lack of optimization robustness against gradient outliers.
- We introduce a novel *algorithmically robust* orthogonalization scheme via an adaptive Newton-Schulz iteration with dimension-specific coefficients, significantly improving orthogonalization accuracy across diverse network architectures.
- We develop an *optimizationally robust* training framework through proximal optimization with soft-thresholding, providing theoretical guarantees for outlier suppression and stable convergence.
- Extensive experiments on LLM pre-training and fine-tuning demonstrate that ROOT achieves superior performance and faster convergence compared to state-of-the-art optimizers, while maintaining computational efficiency, particularly in noisy and non-convex scenarios.

2. Related Work

In this section, we provide a brief overview of related work in the field of optimizers, as well as recent advances.

2.1. Matrix-Aware Optimizers

The training of deep neural networks has historically relied on stochastic first-order methods. Standard optimizers,

such as SGD with momentum (Polyak, 1964) and AdamW (Kingma, 2014; Loshchilov & Hutter, 2017), treat model parameters as independent vectors. While computationally efficient, these coordinate-wise updates overlook the rich structural correlations inherent in weight matrices. To address this, second-order methods like K-FAC (Martens & Grosse, 2015) and Shampoo (Gupta et al., 2018; Anil et al., 2020) leverage Kronecker-factored preconditioners to capture parameter geometry. However, despite their theoretical convergence advantages, these methods often incur prohibitive computational and memory overheads in large-scale settings.

Bridging this gap, the Muon optimizer (Jordan et al., 2024) regulate update geometry through matrix orthogonalization rather than explicit curvature approximation. Specifically, Muon employs a Newton-Schulz iteration (Bernstein & Newhouse, 2024; Higham, 2008; Guo & Higham, 2006) to orthogonalize the momentum matrix, which theoretically corresponds to performing steepest descent under the spectral norm (Li & Hong, 2025; Kovalev, 2025). This approach effectively balances structural awareness with computational efficiency, achieving $O(N)$ complexity similar to first-order methods while promoting coherent updates across parameter matrices.

2.2. Recent Advances in Muon Variants

Following Muon’s success, several variants have emerged to extend its capabilities across different dimensions:

Efficiency and Scalability. To mitigate communication overhead in distributed settings, Dion (Ahn et al., 2025) replaces the dense Newton-Schulz iteration with amortized power iteration. For computational efficiency, LiMuon (Huang et al., 2025) leverages randomized SVD, while Drop-Muon (Gruntkowska et al., 2025) explores randomized layer subsampling to reduce the update frequency.

Adaptivity and Precision. While Muon lacks element-wise adaptivity, recent works attempt to reintegrate it. AdaGO (Zhang et al., 2025) combines orthogonal directions with AdaGrad-style step sizes. AdaMuon (Si et al., 2025) incorporates a second-moment estimator, employing sign-based orthogonalization to enforce stability. On the numerical front, CANS (Grishina et al., 2025) utilizes Chebyshev polynomials to accelerate the convergence of the orthogonalization process over spectral intervals. Complementing these spectral-focused approaches, our work investigates orthogonalization precision across varying matrix dimensions.

3. Approach

This section first outlines the preliminaries and then introduces our method for robust optimization which is achieved via adaptive Newton iteration with fine-grained coefficients

and outlier suppression.

3.1. Preliminaries

The orthogonalization-based optimizers, *e.g.*, Muon (Jordan et al., 2024), address the optimization of neural network parameters that exhibit a matrix structure. During each iteration t , the algorithm updates the weight matrix θ_{t-1} using the current momentum μ , learning rate η_t , and objective function \mathcal{L} . The optimization procedure is defined by the following steps:

$$\begin{aligned} M_t &= \mu M_{t-1} + \nabla \mathcal{L}(\theta_{t-1}) \\ M'_t &= \text{Newton-Schulz}(M_t) \\ \theta_t &= \theta_{t-1} - \eta_t M'_t \end{aligned} \quad (1)$$

In this formulation, M_t represents the gradient momentum at step t , initialized as a zero matrix when $t = 0$. The key innovation lies in the application of a Newton-Schulz (NS) iterative method (Bernstein & Newhouse, 2024) to approximate the transformation $(M_t M_t^T)^{-1/2} M_t$. Considering the singular value decomposition (SVD) of $M_t = U \Sigma V^T$, this transformation yields UV^T , effectively orthogonalizing the momentum matrix. This orthogonalization process promotes isomorphic update matrices, which encourages the network to explore diverse optimization directions rather than converging along a limited set of dominant pathways.

The iterative process begins by initializing $X_0 = M_t / \|M_t\|_F$. At each subsequent step k , the matrix X_k is updated from X_{k-1} according to the recurrence relation:

$$\begin{aligned} X_k &= aX_{k-1} + bX_{k-1}(X_{k-1}^T X_{k-1}) \\ &\quad + cX_{k-1}(X_{k-1}^T X_{k-1})^2 \end{aligned} \quad (2)$$

After N iterations, the resulting matrix X_N serves as the approximation. The coefficients a , b , and c are carefully selected to ensure proper convergence of the iterative scheme. The Muon optimizer adopts the values $a = 3.4445$, $b = -4.7750$, and $c = 2.0315$, which were originally designed to accelerate convergence for matrices with small initial singular values in 5 steps.

Algorithm 1 Muon Optimizer (Jordan et al., 2024)

Require: Learning rate η , momentum μ

- 1: Initialize $M_0 \leftarrow 0$
 - 2: **for** $t = 1, \dots$ **do**
 - 3: Compute gradient $G_t \leftarrow \nabla_{\theta} \mathcal{L}_t(\theta_{t-1})$
 - 4: $M_t \leftarrow \mu M_{t-1} + G_t$
 - 5: $M'_t \leftarrow \text{NewtonSchulz5}(M_t)$
 - 6: Update parameters $\theta_t \leftarrow \theta_{t-1} - \eta M'_t$
 - 7: **end for**
 - 8: **Return** θ_t
-

3.2. Adaptive Newton Iteration with Fine-grained Coefficients

3.2.1. ENHANCING ROBUSTNESS AGAINST MATRIX DIMENSION VARIATIONS

The Muon optimizer’s primary innovation lies in its use of the Newton-Schulz iteration to approximate an orthogonal matrix from the momentum matrix M_t . While computationally efficient, we identify a critical *robustness* limitation: the fixed-coefficient NS iteration exhibits significant sensitivity to matrix dimensions, leading to inconsistent orthogonalization quality across different layers of deep neural networks.

The core issue stems from the one-size-fits-all approach to orthogonalization coefficients. As demonstrated in Table 1, the mean squared error (MSE) of the orthogonal approximation varies dramatically with matrix shape. Square matrices ($n = m$) consistently yield the highest MSE values—up to two orders of magnitude worse than highly non-square configurations. For instance, with $m = 2048, n = 2048$, the MSE drops from 0.0499 to 0.0352 when the coefficients are learned adaptively instead of being fixed. This dimensional sensitivity creates an inherent *fragility* in the optimization process, as layers with different dimensions receive orthogonalization of varying quality, compromising the consistency and reliability of gradient updates.

Table 1. Orthogonalization error reveals dimensional fragility of fixed-coefficient Newton-Schulz iteration.

$n \backslash m$	2048	4096	8192	2048	2048	2048	2048
m	2048	4096	8192	3072	4096	8192	16384
a	3.4445	3.4445	3.4445	3.4445	3.4445	3.4445	3.4445
b	-4.7750	-4.7750	-4.7750	-4.7750	-4.7750	-4.7750	-4.7750
c	2.0315	2.0315	2.0315	2.0315	2.0315	2.0315	2.0315
MSE	0.0499	0.0637	0.0761	0.0338	0.0362	0.0340	0.0425
a	3.3334	3.3732	3.3886	2.9091	2.7739	2.5925	2.7045
b	-4.2591	-4.6134	-4.9026	-3.8108	-3.4911	-3.0010	-3.2335
c	1.7791	2.0576	2.3105	1.8600	1.6992	1.4138	1.5336
MSE	0.0352	0.0470	0.0587	0.0024	0.0010	0.0003	0.0003

To address this dimensional fragility and build a *dimension-robust* orthogonalization process, we propose an adaptive Newton-Schulz iteration (AdaNewton) with fine-grained, dimension-specific coefficients. Instead of using global constants that are optimal only for an “average” matrix shape, we learn specialized coefficients $\{a^{(m,n)}, b^{(m,n)}, c^{(m,n)}\}$ for each unique matrix size (m, n) in the network architecture. This approach ensures consistent high-precision orthogonalization regardless of layer dimensions, making the optimization process robust to the inherent architectural variations in modern neural networks.

Formally, for a matrix $X_{k-1} \in \mathbb{R}^{m \times n}$, our robust adaptive update rule becomes:

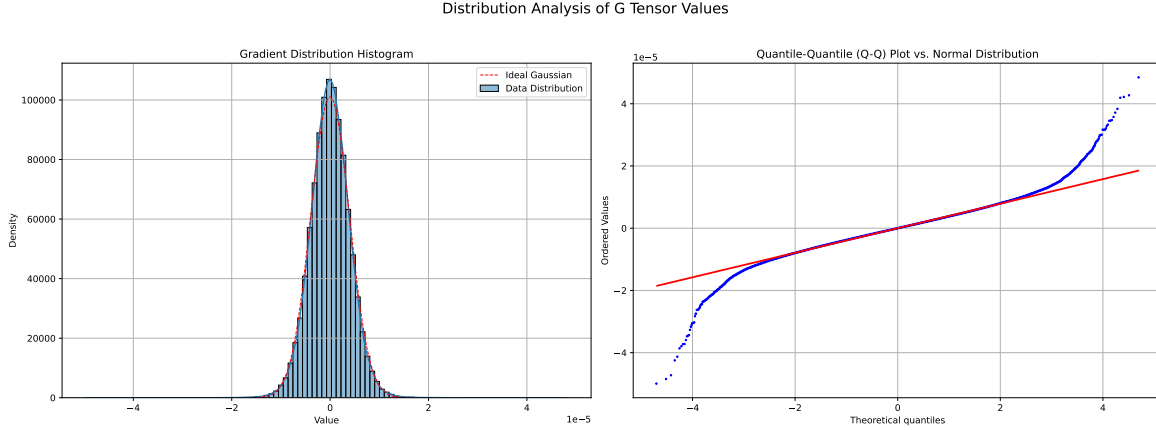


Figure 1. Analysis of gradient distribution revealing outlier characteristics. **(Left)** Histogram with Gaussian reference shows long-tailed distribution. **(Right)** Q-Q plot quantifies deviation from normality, where points deviating from the diagonal indicate outliers. These outliers can disproportionately influence the optimization process.

$$X_k = a^{(m,n)} X_{k-1} + b^{(m,n)} X_{k-1} (X_{k-1}^T X_{k-1}) + c^{(m,n)} X_{k-1} (X_{k-1}^T X_{k-1})^2 \quad (3)$$

The coefficients can be optimized jointly with the model parameters during training, allowing the orthogonalization process to automatically adapt to the specific spectral properties of each layer type. This fine-grained adaptation represents a paradigm shift from fragile dimension-sensitive orthogonalization to robust dimension-invariant orthogonalization, ensuring consistent update quality throughout the network.

3.2.2. THEORETICAL GUARANTEES FOR ROBUST CONVERGENCE

To establish the theoretical superiority of our adaptive coefficient design, we analyze the convergence properties of the Newton-Schulz iteration from a robust optimization perspective. The core insight is that by tailoring coefficients to specific matrix dimensions, we achieve a more precise approximation of the desired orthogonal transformation.

Consider the Newton-Schulz iteration defined by the polynomial mapping:

$$g(x) = ax + bx^3 + cx^5$$

After T iterations, the cumulative transformation is given by the T -fold composition $g^{(T)}(x) = g(g(\dots g(x) \dots))$.

Now, assuming that X_t can be decomposed via SVD as $X_t = U \Sigma_t V^T$, substituting into the iteration yields:

$$X_{t+1} = U(a \Sigma_t + b \Sigma_t^3 + c \Sigma_t^5) V^T$$

Thus, the iteration effectively operates on the singular values through the polynomial $g(x)$. After T iterations, we

have $\Sigma_T = g^{(T)}(\Sigma_0)$, and the optimization objective is to minimize the Frobenius norm distance to the identity matrix:

$$\mathcal{L}_{\text{newton}}(a, b, c) = \|g^{(T)}(\Sigma) - I\|_F^2 \quad (4)$$

The standard fixed-coefficient approach selects parameters (a, b, c) that minimize the worst-case error over a global singular value interval $I_{\text{std}} = [\sigma_{\min}, \sigma_{\max}]$:

$$(a^*, b^*, c^*) = \arg \min_{a, b, c} \max_{\sigma \in I_{\text{std}}} |g^{(T)}(\sigma) - 1|$$

This yields a convergence guarantee with error bound E_{std} .

However, this one-size-fits-all approach is inherently fragile because real-world weight matrices of different dimensions exhibit distinct singular value distributions $S^{(m,n)} \subseteq I_{\text{std}}$. Our robust adaptive method learns dimension-specific coefficients that minimize the error over these characteristic distributions:

$$(a^{(m,n)}, b^{(m,n)}, c^{(m,n)}) = \arg \min_{a, b, c} \max_{\sigma \in S^{(m,n)}} |g^{(T)}(\sigma) - 1|$$

The key theoretical advantage emerges from the optimization structure. Since $S^{(m,n)} \subseteq I_{\text{std}}$, the minimax error over the smaller set is bounded by the global error:

$$\begin{aligned} E_{(m,n)} &= \min_{a, b, c} \max_{\sigma \in S^{(m,n)}} |g^{(T)}(\sigma) - 1| \\ &\leq \min_{a, b, c} \max_{\sigma \in I_{\text{std}}} |g^{(T)}(\sigma) - 1| = E_{\text{std}} \end{aligned}$$

Moreover, when $S^{(m,n)}$ is a proper subset of I_{std} (which occurs for non-square matrices), the adaptive coefficients can achieve strictly better approximation. Formally, if $S^{(m,n)} \subset I_{\text{std}}$ and the function $g^{(T)}(x)$ is not constant, then typically:

$$E_{(m,n)} < E_{\text{std}}$$

This improvement translates directly to our optimization objective $\|g^{(T)}(\Sigma) - I\|_F^2$. For a matrix with singular values $\{\sigma_i\}$, we have:

$$\|g^{(T)}(\Sigma) - I\|_F^2 = \sum_i (g^{(T)}(\sigma_i) - 1)^2$$

By achieving a smaller maximum error $\max_i |g^{(T)}(\sigma_i) - 1|$ through our adaptive coefficients, we obtain a tighter bound on the Frobenius norm objective in Eq. 4. This theoretical result demonstrates that our dimension-adaptive approach provides provably better orthogonalization compared to the fixed-coefficient method, establishing the robustness advantages of our design.

3.3. Robust Optimization via Outlier Suppression

3.3.1. ENHANCING ROBUSTNESS AGAINST OUTLIER-INDUCED GRADIENT NOISE

In large-scale language model training, mini-batch gradients are frequently contaminated by outlier noise with gradient components of anomalously large magnitudes (Zhang et al., 2025). Figure 1 provides a conceptual visualization of such outlier noise within a gradient distribution.

These outliers pose a critical threat to the stability of the orthogonalization process in Muon. The NS algorithm requires an initial normalization of the input matrix M_t to ensure its singular values are within a specific range for convergence. The presence of extreme outlier noise in M_t can distort this normalization step and, more critically, be amplified by the polynomial nature of the NS iteration. This amplification compromises the intended denoising effect of orthogonalization, as erratic directions associated with outliers are preserved in the parameter update, potentially destabilizing training and causing issues like exploding attention logits in Transformers.

To build an *outlier-robust* optimization framework, we incorporate proximal optimization with soft-thresholding. This approach provides inherent protection against gradient contamination while maintaining the benefits of momentum orthogonalization.

3.3.2. SOFT-THRESHOLDING FOR OUTLIER SUPPRESSION

Our robust optimization framework is derived from a gradient clipping perspective that explicitly handles outlier noise while preserving the overall gradient direction. We model the gradient (or momentum) M_t as comprising two components: a base component B_t containing the well-behaved gradient information, and an outlier component O_t representing anomalous large-magnitude elements:

$$M_t = B_t + O_t \quad (5)$$

The robust optimization objective aims to separate these components while constraining the influence of outliers:

$$\min_{B_t, O_t} \|M_t - B_t - O_t\|_F^2 + \lambda \|O_t\|_1 \quad \text{subject to} \quad \|B_t\| \leq \tau \quad (6)$$

where λ controls the sparsity of outlier detection and τ bounds the magnitude of the base component. This formulation explicitly penalizes the presence of large outlier components while ensuring the base component remains within reasonable bounds.

The optimal solution to this robust decomposition is characterized by the proximal operator for the L1-norm, which yields the soft-thresholding function (Tibshirani, 1996; Donoho, 2002):

$$\mathcal{T}_\varepsilon[x]_i = \text{sign}(x_i) \cdot \max(|x_i| - \varepsilon, 0) \quad (7)$$

where ε is the threshold hyperparameter related to the regularization parameter λ . The soft-thresholding can be expressed as

$$\mathcal{T}_\varepsilon[x]_i = \begin{cases} x_i - \varepsilon, & \text{if } x_i > \varepsilon \\ 0, & \text{if } |x_i| \leq \varepsilon \\ x_i + \varepsilon, & \text{if } x_i < -\varepsilon \end{cases} \quad (8)$$

This operation provides the closed-form solution for outlier separation:

$$\begin{cases} O_t = \mathcal{T}_\varepsilon(M_t), & (\text{sparse outlier components}) \\ B_t = M_t - O_t. & (\text{clipped, robust components}) \end{cases} \quad (9)$$

Mathematically, soft-thresholding can be interpreted as a continuous, differentiable alternative to hard clipping. While traditional gradient clipping abruptly truncates values beyond a fixed threshold, soft-thresholding applies a smooth shrinkage operation that preserves the relative ordering of gradient magnitudes while dampening extreme values.

Algorithm 2 ROOT Optimizer

Require: Learning rate η , momentum μ , threshold ε

- 1: Initialize $M_0 \leftarrow 0$
 - 2: **for** $t = 1, \dots$ **do**
 - 3: Compute gradient $G_t \leftarrow \nabla_\theta \mathcal{L}(\theta_{t-1})$
 - 4: $M_t \leftarrow \mu M_{t-1} + G_t$ # Momentum accumulation
 - 5: $O_t \leftarrow \mathcal{T}_\varepsilon[M_t]$ # Outlier separation via soft-thresholding
 - 6: $B_t \leftarrow M_t - O_t$ # Clipped base components
 - 7: $B_t^{\text{orth}} \leftarrow \text{AdaNewton}(B_t)$ # Robust orthogonalization
 - 8: Update parameters $\theta_t \leftarrow \theta_{t-1} - \eta B_t^{\text{orth}}$
 - 9: **end for**
 - 10: **Return** θ_t
-

In the context of ROOT, we apply this robust decomposition to the momentum matrix $M_t \in \mathbb{R}^{m \times n}$:

$$[\mathcal{T}_\varepsilon(M_t)]_{ij} = \mathcal{T}_\varepsilon([M_t]_{ij}) \quad (10)$$

The key innovation is that orthogonalization is applied only to the robust component B_t , while outlier components O_t are discarded. This ensures that the orthogonalization process which is highly sensitive to large magnitude variations, operates on stable and clipped gradients, dramatically improving training robustness, as shown in Algorithm 2.

4. Experiments

4.1. Implementation Details

Experiments Setup To validate our approach, we conduct pretraining using the FineWeb-Edu dataset (Penedo et al., 2024), utilizing a 10-billion-token subset for ablation studies and a 100-billion-token sample for the main experiment. Both subsets are available at (HuggingFaceFW, 2024). To systematically evaluate the proposed optimization enhancements, we train a 1B Transformer (Dubey et al., 2024; Rang et al., 2025). All models are trained on distributed clusters of Ascend NPUs. The training infrastructure leverages high-speed interconnects to enable efficient data and model parallelism, ensuring scalable training across multiple nodes. An attention-mask-reset strategy (Chen et al., 2025) is also used to prevent self-attention between different documents within a sequence.

Regarding the hyperparameters, all models are pretrained for a single epoch. We employ a cosine learning rate schedule that decays to 10% of the peak learning-rate, following a warm-up phase of 2,000 steps. Specifically, for the 10B-token ablation experiments, the peak learning rate is set to 8×10^{-4} with a global batch size of 0.4M. The pre-training sequence length is set to 4096. For the 100B-token main experiments, the learning rate is increased to 1.6×10^{-3} , and the global batch size is set to 1M. We adopt the default Muon hyperparameters from (Jordan et al., 2024) and apply a 0.2 scaling factor following (Liu et al., 2025) to align the update RMS with AdamW.

Evaluation Benchmark Setup We evaluate our method on a comprehensive set of Academic benchmarks: Hel-laSwag (Zellers et al., 2019), ARC-easy (ARC-e) and ARC-challenge (ARC-c) (Clark et al., 2018), BoolQ (Clark et al., 2019), PIQA (Bisk et al., 2020), SciQ (Welbl et al., 2017), WINO (Sakaguchi et al., 2019), OBQA (Mihaylov et al., 2018), and WSC (Levesque et al., 2012). We utilize the lm-evaluation-harness framework (Gao et al., 2024) for all evaluations, and we evaluate all tasks in zero-shot.

4.2. Validation on Real Training Dynamics

While Table 1 establishes the benefits of shape-specific coefficients on static distributions, we further validate whether these coefficients generalize to the dynamic spectral shifts observed during actual LLM training. We sampled input gradients G_t from the first 10k pre-training steps, covering phases from early instability to stable convergence.

In this experiment, we compare three distinct orthogonalization strategies, each restricted to a strict budget of 5 iterations. We evaluate the Muon Baseline (Jordan et al., 2024) with its general-purpose coefficients ($a = 3.4445, b = -4.7750, c = 2.0315$), as well as the Classic Newton-Schulz (Quintic) which uses coefficients (1.875, -1.25, 0.375) for order-5 convergence. We contrast these fixed strategies with our proposed ROOT (AdaNewton), which utilizes shape-specific coefficients learned from training gradients to explicitly target the geometry of different layers.

To quantify approximation fidelity, we report the Relative Error ($\|\hat{O} - O_{\text{SVD}}\|_F / \|O_{\text{SVD}}\|_F$) relative to the ground-truth SVD. The reported error is averaged across all parameter types targeted by the optimizer (Attention Query/Output and MLP Up/Down projections) at each step.

Figure 2 illustrates the evolution of this metric. ROOT (Red) consistently maintains a lower relative error compared to the baselines throughout the training trajectory. While the Muon baseline (Blue) exhibits a distinguishable error floor, likely due to the sub-optimality of fixed coefficients for specific aspect ratios, ROOT achieves a closer approximation to the ground truth. These results suggest that incorporating dimension-aware coefficients mitigates the precision loss often associated with fixed-coefficient Newton-Schulz iterations on diverse matrix shapes.

4.3. LLM Pre-training Analysis

We evaluate the pre-training performance of the proposed ROOT optimizer on the 1B Transformer over a 10B-token trajectory. Figure 3 visualizes the training loss. We compare the Muon baseline against two configurations: *ROOT-SoftThresh* (outlier suppression only) and the full ROOT optimizer.

As observed, both ROOT variants achieve consistently lower training loss compared to the Muon baseline. The ROOT (SoftThresh) variant improves convergence over standard Muon, suggesting that suppressing gradient outliers aids in stabilizing the optimization process. Furthermore, the full ROOT optimizer yields the lowest loss, indicating that dimension-adaptive orthogonalization provides additive benefits to outlier suppression. Ultimately, ROOT reaches a final training loss of 2.5407, surpassing the Muon baseline by 0.01.

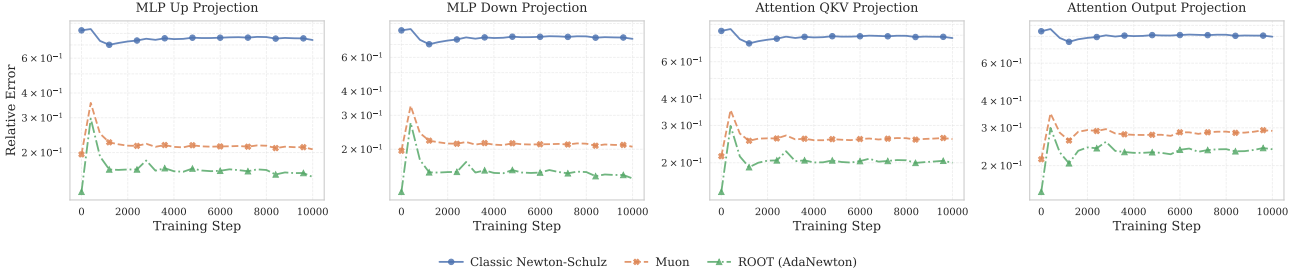


Figure 2. Orthogonalization precision relative to ground-truth SVD. The plot tracks the Relative Error averaged over all optimized parameters (Attention QKV/O and MLP Up/Down projections). Under a fixed 5 iterations, ROOT maintains lower approximation error compared to the Muon baseline and Classic Newton-Schulz. This indicates that shape-specific coefficients provide superior fidelity across varying matrix dimensions.

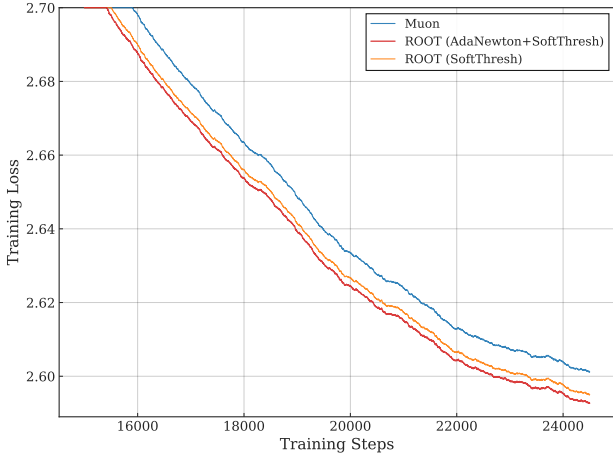


Figure 3. Training loss curves for 10B tokens. ROOT variants demonstrate faster convergence and lower final loss compared to Muon baseline, with full ROOT achieving the best performance.

4.4. Benchmark Evaluation Results

To assess model generalization, we further train the 1B model on 100B tokens and evaluate the trained 1B models across a diverse set of common academic benchmarks. Table 2 summarizes the zero-shot performance. ROOT achieves competitive or superior performance across the evaluated tasks. These results confirm that ROOT enhancements not only accelerate training convergence but also yield higher-quality final language models.

4.5. Ablation Studies

Outlier Suppression Threshold. Instead of a fixed threshold scalar, we employ a dynamic, percentile-based threshold $\varepsilon_t = \text{Quantile}(|M_t|, p)$ to adapt to the evolving gradient scale. We investigate the sensitivity of convergence to varying percentiles $p \in \{0.85, 0.90, 0.95, 0.99\}$. As illustrated

in Figure 4, the results reveal a trade-off between outlier suppression and signal preservation. A conservative threshold ($p = 0.99$) results in under-suppression, failing to filter the heavy tail of the noise distribution and leading to sub-optimal stability. Conversely, an overly aggressive threshold ($p = 0.85$) causes excessive truncation. We identify $p = 0.90$ as the optimal equilibrium for LLM pre-training tasks, effectively isolating outliers while retaining the structural integrity of the gradients.

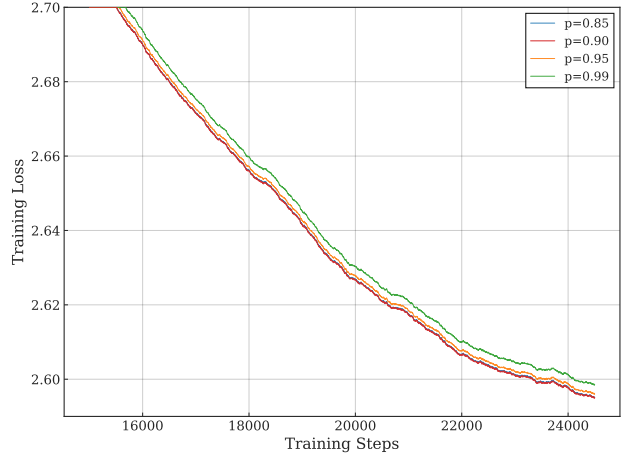


Figure 4. Ablation on the quantile hyperparameter p . The curve with $p = 0.90$ demonstrates the optimal equilibrium between suppressing gradient noise and preserving informative gradient signals.

Spectral Calibration Strategy. To determine the optimal coefficients $\{a, b, c\}$ for ROOT (AdaNewton), we perform offline optimization using singular value distributions of momentum matrices with varying dimensions processed by the Muon optimizer, collected during the training process. We evaluate three calibration Strategies: Random Calibration, optimized solely on random Gaussian matrices; Mixed

Table 2. Zero-shot performance on standard LLM benchmarks. ROOT outperforms both the AdamW baseline and the Muon optimizer across diverse common academic tasks.

METHOD	HELLASWAG	BOOLQ	PIQA	ARC-E	ARC-C	OBQA	SCIQ	WINO	WSC	AVG.
ADAMW	44.24	62.60	72.69	71.63	37.80	27.20	89.80	58.09	67.40	59.05
MUON	44.83	61.16	73.07	74.12	37.12	29.80	89.50	59.67	67.03	59.59
ROOT	45.37	62.08	73.12	72.14	36.86	31.20	90.40	60.30	69.60	60.12

Calibration (1:1), where real `ns_input` samples are augmented with random matrices at an equal ratio; and Mixed Calibration (1:3), which increases the proportion of random matrices to a 3:1 ratio.

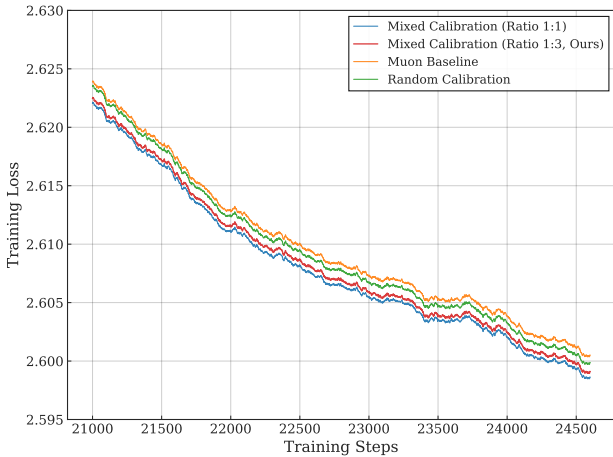


Figure 5. Ablation on the data composition for coefficient calibration. While a high ratio of real samples (Ratio 1:1) achieves lower loss here, it induces instability (loss spikes) in larger-scale experiments or when combined with ROOT (SoftThresh). The Mixed (1:3) strategy provides the optimal balance between convergence speed and robustness.

Figure 5 illustrates the training trajectories under these configurations. While the Mixed (1:1) strategy yields the lowest loss in this specific setting, broader evaluations reveal its vulnerability to instability (e.g., loss spikes) when scaled to larger models or integrated with the Soft-Thresholding mechanism. Conversely, the Random Calibration baseline offers limited convergence acceleration. Consequently, we adopt the Mixed (1:3) strategy, which effectively prevents overfitting, thereby balancing accelerated convergence with generalization stability.

4.6. Generalization to Vision Tasks

To evaluate the generalization capabilities of ROOT beyond language modeling, we conducted image classification experiments by training a lightweight Vision Transformer ($\approx 6.3\text{M}$ parameters) on CIFAR-10 from scratch.

We adopt a compact architecture adapted from (Omihub777, 2023), which processes images via 4×4 patches. In this experiment, we explicitly isolate the efficacy of the Soft-Thresholding mechanism against the Muon baseline. Models were trained for 100 epochs. Following standard protocols, all 2D weight matrices were optimized via Muon or ROOT, while 1D parameters (biases, norms) and embeddings and class-tokens were optimized using AdamW. The results in Table 3 show that ROOT consistently outperforms the baseline. The improvement is most significant for the quantile percentile hyperparameter of 0.85, where the baseline achieves only 84.67% accuracy. These results confirm that the soft-thresholding mechanism mitigates gradient noise, thereby enhancing generalization even in non-language modalities.

Table 3. Top-1 Test Accuracy on CIFAR-10 (ViT, 6.3 M, trained from scratch). Fixed params: Muon LR = 0.02, AdamW LR = 0.001, WD = 5×10^{-5} .

METHOD	p (QUANTILE)	ACC (%)
MUON	-	84.67
ROOT	0.95	85.75
ROOT	0.90	86.58
ROOT	0.85	88.44

5. Conclusion

In this work, we have presented ROOT, a robust orthogonalized optimizer that addresses two critical limitations in modern large-scale language model training. By introducing dimension-robust orthogonalization through adaptive Newton iterations and optimization robustness via proximal outlier suppression, ROOT establishes a new paradigm for stable and efficient neural network optimization. Our extensive experimental validation demonstrates that the proposed method achieves superior performance in challenging noisy and non-convex scenarios, providing both theoretical guarantees and practical benefits. This work opens promising directions for developing robust optimization frameworks that can handle the increasing complexity and scale of future language models, potentially enabling more reliable and efficient training of next-generation AI systems.

References

- Ahn, K., Xu, B., Abreu, N., Fan, Y., Magakyan, G., Sharma, P., Zhan, Z., and Langford, J. Dion: Distributed orthonormalized updates. *arXiv preprint arXiv:2504.05295*, 2025.
- Anil, R., Gupta, V., Koren, T., Regan, K., and Singer, Y. Scalable second order optimization for deep learning. *arXiv preprint arXiv:2002.09018*, 2020.
- Bernstein, J. and Newhouse, L. Old optimizer, new norm: An anthology. *arXiv preprint arXiv:2409.20325*, 2024.
- Bisk, Y., Zellers, R., Bras, R. L., Gao, J., and Choi, Y. Piqa: Reasoning about physical commonsense in natural language. In *Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901, 2020.
- Chen, H., Wang, Y., Han, K., Li, D., Li, L., Bi, Z., Li, J., Wang, H., Mi, F., Zhu, M., et al. Pangu embedded: An efficient dual-system llm reasoner with metacognition. *arXiv preprint arXiv:2505.22375*, 2025.
- Clark, C., Lee, K., Chang, M.-W., Kwiakowski, T., Collins, M., and Toutanova, K. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In Burstein, J., Doran, C., and Solorio, T. (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 2924–2936, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1300. URL <https://aclanthology.org/N19-1300/>.
- Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., and Tafjord, O. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.
- Donoho, D. L. De-noising by soft-thresholding. *IEEE transactions on information theory*, 41(3):613–627, 2002.
- 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.
- Gao, L., Tow, J., Abbasi, B., Biderman, S., Black, S., DiPofi, A., Foster, C., Golding, L., Hsu, J., Le Noac’h, A., Li, H., McDonnell, K., Muennighoff, N., Ociepa, C., Phang, J., Reynolds, L., Schoelkopf, H., Skowron, A., Sutawika, L., Tang, E., Thite, A., Wang, B., Wang, K., and Zou, A. The language model evaluation harness, 07 2024. URL <https://zenodo.org/records/12608602>.
- Grishina, E., Smirnov, M., and Rakhuba, M. Accelerating newton-schulz iteration for orthogonalization via chebyshev-type polynomials. *arXiv preprint arXiv:2506.10935*, 2025.
- Grutkowska, K., Maziane, Y., Qu, Z., and Richtárik, P. Drop-muon: Update less, converge faster. *arXiv preprint arXiv:2510.02239*, 2025.
- Guo, C.-H. and Higham, N. J. A schur-newton method for the matrix p th root and its inverse. *SIAM Journal on Matrix Analysis and Applications*, 28(3):788–804, 2006.
- Gupta, V., Koren, T., and Singer, Y. Shampoo: Pre-conditioned stochastic tensor optimization. In *International Conference on Machine Learning*, pp. 1842–1850. PMLR, 2018.
- Higham, N. J. *Functions of matrices: theory and computation*. SIAM, 2008.
- Huang, F., Luo, Y., and Chen, S. Limuon: Light and fast muon optimizer for large models. *arXiv preprint arXiv:2509.14562*, 2025.
- HuggingFaceFW. fineweb-edu (revision 22b0aca), 2024. URL <https://huggingface.co/datasets/HuggingFaceFW/fineweb-edu>.
- 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/>.
- Kingma, D. P. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Kovalev, D. Understanding gradient orthogonalization for deep learning via non-euclidean trust-region optimization. *arXiv preprint arXiv:2503.12645*, 2025.
- Levesque, H., Davis, E., and Morgenstern, L. The winograd schema challenge. In *13th International Conference on the Principles of Knowledge Representation and Reasoning, KR 2012*, Proceedings of the International Conference on Knowledge Representation and Reasoning, pp. 552–561. Institute of Electrical and Electronics Engineers Inc., 2012. ISBN 9781577355601. 13th International Conference on the Principles of Knowledge Representation and Reasoning, KR 2012 ; Conference date: 10-06-2012 Through 14-06-2012.
- Li, J. and Hong, M. A note on the convergence of muon and further. *arXiv e-prints*, pp. arXiv–2502, 2025.

- Liu, J., Su, J., Yao, X., Jiang, Z., Lai, G., Du, Y., Qin, Y., Xu, W., Lu, E., Yan, J., et al. Muon is scalable for llm training. *arXiv preprint arXiv:2502.16982*, 2025.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Martens, J. and Grosse, R. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pp. 2408–2417. PMLR, 2015.
- Mihaylov, T., Clark, P., Khot, T., and Sabharwal, A. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *EMNLP*, 2018.
- Omihub777. ViT-CIFAR: Vision transformer for CIFAR-10/100 on pytorch. GitHub repository, 2023. URL <https://github.com/omihub777/ViT-CIFAR>. Commit main branch, 2023-11-20.
- OpenAI. Gpt-4 technical report. Technical report, OpenAI, 2023. URL <https://cdn.openai.com/papers/gpt-4.pdf>. Accessed: 2024-11-04.
- Penedo, G., Kydlíček, H., Lozhkov, A., Mitchell, M., Raffel, C. A., Von Werra, L., Wolf, T., et al. The fineweb datasets: Decanting the web for the finest text data at scale. *Advances in Neural Information Processing Systems*, 37: 30811–30849, 2024.
- Polyak, B. T. Some methods of speeding up the convergence of iteration methods. *Ussr computational mathematics and mathematical physics*, 4(5):1–17, 1964.
- Rang, M., Bi, Z., Zhou, H., Chen, H., Xiao, A., Guo, T., Han, K., Chen, X., and Wang, Y. Revealing the power of post-training for small language models via knowledge distillation. *arXiv preprint arXiv:2509.26497*, 2025.
- Robbins, H. and Monro, S. A stochastic approximation method. *The annals of mathematical statistics*, pp. 400–407, 1951.
- Sakaguchi, K., Bras, R. L., Bhagavatula, C., and Choi, Y. Winogrande: An adversarial winograd schema challenge at scale. *arXiv preprint arXiv:1907.10641*, 2019.
- Si, C., Zhang, D., and Shen, W. Adamuon: Adaptive muon optimizer. *arXiv preprint arXiv:2507.11005*, 2025.
- Tibshirani, R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 58(1):267–288, 1996.
- Welbl, J., Liu, N. F., and Gardner, M. Crowdsourcing multiple choice science questions. In *NUT@EMNLP*, 2017.
- Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., and Choi, Y. HellaSwag: Can a machine really finish your sentence? In Korhonen, A., Traum, D., and Màrquez, L. s. (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4791–4800, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1472. URL <https://aclanthology.org/P19-1472/>.
- Zhang, M., Liu, Y., and Schaeffer, H. Adagrad meets muon: Adaptive stepsizes for orthogonal updates. *arXiv preprint arXiv:2509.02981*, 2025.