

Comprehensive Analysis of Optimization Algorithms

Jinghao Liu (jliu63) Xuan Zhang (xuanz24)
Yuzheng Zhang (yuzhez4)

December 2025

Abstract

We conduct a unified empirical and theoretical comparison of modern deep-learning optimizers, focusing on Adam-based methods and the recently proposed Muon. Using CIFAR-10/100 under controlled hyperparameters, consistent augmentation, and multiple seeds, we evaluate SGD with momentum, Adam, AdamW, RAdam, Lion, and Muon. Beyond final accuracy, we study convergence, stability, generalization, and failure modes, supported by four targeted experiments on stability (H1), regularization (H2), label-noise robustness (H3), and data efficiency (H4).

Our findings: (1) well-tuned SGD with momentum still achieves the highest accuracy; (2) AdamW reliably outperforms Adam and should be the default adaptive optimizer; (3) RAdam improves early-stage stability but is sensitive to implementation details like weight decay; and (4) Muon performs strongly in low-data settings and is competitive on CIFAR-100. We conclude with practical optimizer guidelines and directions such as dynamic optimizer switching.

1 Introduction

While adaptive optimizers like Adam [1], Lion [4], and the geometry-aware Muon [5] offer theoretical gains over SGD, empirical adoption often defaults to Adam without rigorous comparison. Furthermore, implementation nuances create a significant gap between theoretical claims and practical performance. We address this by establishing a unified vision benchmark characterized by:

- **Strict Control:** Identical architectures and schedules to ensure fair comparison.
- **Statistical Rigor:** Aggregated metrics across multiple random seeds.
- **Hypothesis Verification:** Targeted experiments probing specific theoretical claims.

2 Related Work

Adaptive gradient methods. Adam [1] introduced exponential moving averages of first and second moments, enabling per-parameter adaptive learning rates that improve stability and speed. Its coupled L2 regularization and generalization issues motivated AdamW [2], which decouples weight decay from the adaptive update. RAdam [3] analyzes the variance

of adaptive learning rates and adds a rectification term to stabilize early training when second-moment estimates are unreliable.

Sign-based and geometry-aware optimizers. Lion [4] uses the sign of a momentum-like term, reducing memory use and potentially increasing robustness to outliers. Muon [5, 6] takes a geometry-aware approach for hidden weight matrices, using a Newton–Schulz iteration to orthogonalize updates and normalize them in spectral norm. It is typically paired with AdamW for non-matrix parameters (e.g., embeddings and heads).

Optimizer benchmarks. Existing comparisons often cover only a subset of optimizers, a single dataset, or inconsistent hyperparameters. Our work differs by (i) evaluating a broad suite including Muon, (ii) using unified and controlled protocols, and (iii) conducting hypothesis-driven tests grounded in theory.

3 Experimental Setup

3.1 Datasets

We use two standard image classification benchmarks:

- **CIFAR-10:** 50,000 training and 10,000 test images of size 32×32 with 10 classes. This dataset is relatively easy; most modern architectures achieve above 90% accuracy.
- **CIFAR-100:** 50,000 training and 10,000 test images with 100 fine-grained classes. Each class has only 500 training examples, making it more sensitive to regularization, overfitting, and optimizer behavior.

Both datasets are loaded using `torchvision.datasets`, with automatic downloading and caching.

3.2 Models

To match typical practice, we use different architectures for the two datasets:

- **ResNet-18:** For CIFAR-10, an 18-layer residual network with ~ 11 M parameters.
- **WideResNet-16-4:** for CIFAR-100, a 16-layer residual network with width factor 4.

For CIFAR-100, we enable dropout in WRN-16-4 to further stress-test the regularization behavior of optimizers.

3.3 Data Augmentation and Preprocessing

Unless otherwise stated, we apply the following standard augmentations:

- Random crop with 4-pixel padding.
- Random horizontal flip with probability 0.5.
- Per-channel mean and standard deviation normalization.

For the optimized CIFAR-100 experiments we also explore stronger augmentation (AutoAugment, Cutout) and label smoothing of 0.1; when comparing optimizers within the same setting, all augmentations are kept fixed.

3.4 Optimizers and Hyperparameters

We compare six optimizers:

- SGD with momentum (0.9) and Nesterov disabled.
- Adam with default betas $(\beta_1, \beta_2) = (0.9, 0.999)$.
- AdamW with the same betas but decoupled weight decay.
- RAdam using `torch.optim.RAdam`.
- Lion via the `lion-pytorch` implementation.
- Muon via a custom implementation adapted from the official repository [6], applied to hidden weight matrices and combined with AdamW for other parameters.

Table 1 summarizes the default hyperparameters we use unless otherwise stated.

Optimizer	Base LR	Weight Decay	Notes	Scheduler
SGD	0.10	5×10^{-4}	Momentum 0.9	Cosine
Adam	0.001	0.0	Default betas	Cosine
AdamW	0.001	0.01	Decoupled WD	Cosine
RAdam	0.001	0.001/0.01	See Sec. 4.3	Cosine
Lion	0.0001	0.01	As in [4]	Cosine
Muon	0.02 (hidden)	0.01	Aux AdamW LR 3×10^{-4}	Cosine

Table 1: Default hyperparameters for the six optimizers in our benchmark.

3.5 Training Protocol

For each pair, we run three experiments with different random seeds (42, 123, 456):

- We train for 100 epochs on both CIFAR-10 and CIFAR-100.
- Batch size is 128.
- We use a cosine learning rate scheduler with warmup disabled.
- We log metrics at every epoch: training loss/accuracy, test loss/accuracy, learning rate, and epoch wall-clock time.

All experiments are run on a mix of RTX 3060 and A40.

4 Baseline Results

4.1 CIFAR-10 and CIFAR-100 Summary

Table 2 summarizes performance trends for all optimizers, combining evidence from CIFAR-10 and CIFAR-100.

Several trends emerge consistently across both datasets:

- SGD with momentum attains the highest final accuracy on both CIFAR-10 and CIFAR-100, though its early convergence is slower than adaptive methods.
- AdamW consistently outperforms Adam in both accuracy and generalization, validating the benefit of decoupled weight decay.

Optimizer	C10 Acc.	C100 Acc.	Convergence	Notes
SGD + Momentum	87.85 \pm 0.35	76.91 \pm 0.06	Slow start	Best final accuracy
Muon	87.67 \pm 0.12	75.87 \pm 0.30	Very Fast	Close to SGD on C100
AdamW	86.44 \pm 0.25	74.13 \pm 0.08	Fast	Better generalization
Adam	86.08 \pm 0.09	73.84 \pm 0.18	Fast	Good baseline, worse late
Lion	85.85 \pm 0.05	72.68 \pm 0.06	Fast	Most stable, lower peak
RAdam (v2)	84.76 \pm 0.49	73.21 \pm 0.10	Medium	Early stable, sensitive WD

Table 2: Optimizer performance on CIFAR-10 (ResNet-18) and CIFAR-100 (WRN-16-4), averaged over 3 seeds (42, 123, 456).

- Lion shows the lowest seed variance and strong stability, but sacrifices a small amount of final accuracy.
- RAdam is sensitive to weight-decay implementation; once corrected (v2), it is competitive but still behind AdamW and Muon.
- Muon matches or slightly surpasses SGD on CIFAR-10 and is the strongest adaptive competitor on CIFAR-100, suggesting geometry-aware updates help on more challenging tasks.

These trends are also visible in the learning-curve plots in Appendix 8, particularly Figures 1 and 2.

4.2 Efficiency and Convergence Dynamics

Beyond final accuracy, we analyzed the efficiency of each optimizer. Since experiments were run on heterogeneous hardware (A40, RTX 3060, MPS), we compare ****epochs to target accuracy**** rather than wall-clock time to ensure a fair, hardware-agnostic comparison of sample efficiency. Table 3 highlights a striking advantage for Muon.

Optimizer	C10 Epochs-to-85%	C100 Epochs-to-60%	Early Acc (20ep) C10 / C100
Muon	18.7	9.7	81.05% / 57.39%
AdamW	47.3	18.7	74.95% / 44.85%
SGD	78.3	49.0	66.56% / 36.37%
Lion	65.0	21.3	75.93% / 42.89%

Table 3: Efficiency metrics based on training epochs. Muon reaches functional accuracy thresholds using 50-80% fewer epochs than baselines and learns significantly faster in early stages.

Key Findings:

- **Muon’s Sample Efficiency:** On CIFAR-100, Muon reaches 60% accuracy in just **9.7 epochs**, compared to 18.7 for AdamW and 49.0 for SGD. This represents a **2 \times speedup** over AdamW and **5 \times speedup** over SGD in terms of data passes.
- **Early Learning:** In the first 20 epochs, Muon achieves an average test accuracy of 57.39% on CIFAR-100, nearly 13 points higher than AdamW and 21 points higher than SGD.

- **Overfitting Resistance:** We also monitored the increase in train-test gap from the best epoch to the final epoch. SGD (0.11%) and Muon (0.06%) showed almost no degradation, whereas Adam/AdamW degraded by $\sim 0.36\%$, suggesting that geometry-aware updates provide inherent regularization against late-stage overfitting.

4.3 RAdam Implementation Adjustment

Initial experiments with PyTorch’s `optim.RAdam` yielded anomalously low accuracy ($29.55\% \pm 0.35$) on CIFAR-100. We attribute this to the implementation’s use of coupled L2 regularization rather than decoupled weight decay. At our default `weight_decay` of 10^{-2} , this mechanism severely over-regularized the model relative to the adaptive learning rates. Reducing the decay to 10^{-3} resolved the issue, recovering accuracy to $73.21\% \pm 0.08$. We denote this corrected configuration as *RAdam-v2*.

5 Hypothesis-Driven Experiments

Beyond aggregate benchmarking, we designed four targeted experiments to test specific theoretical claims. Each hypothesis was formulated *a priori*, and we report results regardless of outcome. All results and visualizations are in `results_hypothesis/` and `analysis_hypothesis/`.

5.1 H1: RAdam Early-Stage Stability

Hypothesis. RAdam’s variance rectification will reduce the coefficient of variation (CV) of training loss by $\geq 20\%$ during the first 10 epochs compared to Adam on CIFAR-100.

Theoretical Motivation. Adam’s adaptive learning rate $\alpha_t = \eta/\sqrt{v_t}$ has high variance when t is small. RAdam rectifies this using the degrees of freedom $\rho_t = \rho_\infty - 2t\beta_2^t/(1 - \beta_2^t)$, where $\rho_\infty = 2/(1 - \beta_2)$:

$$r_t = \sqrt{\frac{(\rho_t - 4)(\rho_t - 2)\rho_\infty}{(\rho_\infty - 4)(\rho_\infty - 2)\rho_t}} \quad (1)$$

When $\rho_t < 4$, the variance is effectively infinite; r_t implements automatic warmup.

Setup and Results. We trained WRN-16-4 on CIFAR-100 with matched hyperparameters across three seeds. RAdam-v2 achieved CV of 0.1207 versus Adam’s 0.1342, a **10.1% reduction**. Figure 3 shows visibly smoother loss trajectories for RAdam-v2.

Conclusion. H1 Partially Supported. While below the 20% threshold, the empirical reduction validates RAdam’s core mechanism. The automatic warmup effect is real and measurable, though moderated by batch normalization and residual connections in modern architectures.

5.2 H2: AdamW Regularization Advantage

Hypothesis. AdamW will maintain stable performance under high weight decay ($\lambda = 10^{-2}$), while Adam will suffer catastrophic failure due to coupled L2 regularization.

Theoretical Motivation. Adam with L2 penalty creates per-parameter regularization $\eta\lambda/\sqrt{v_t}$: parameters with large v_t receive *weaker* penalties, while those with small v_t are over-regularized. Standard Adam:

$$\theta_t \leftarrow \theta_{t-1} - \eta \frac{m_t}{\sqrt{v_t} + \epsilon}, \quad m_t \propto g_t + \lambda \theta_{t-1} \quad (2)$$

AdamW decouples weight decay:

$$\theta_t \leftarrow (1 - \eta\lambda)\theta_{t-1} - \eta \frac{m_t}{\sqrt{v_t} + \epsilon} \quad (3)$$

ensuring uniform regularization $\eta\lambda$ across all parameters.

Setup and Results. We trained WRN-16-4 on CIFAR-100 with `weight_decay=0.01` using strong augmentation. **Adam completely collapsed:** 18.67% ± 0.76 (seeds: 17.74%, 19.13%, 19.14%), barely above random (1%). **AdamW remained stable:** 69.68% ± 0.19 (seeds: 69.51%, 69.89%, 69.63%). The **51-point gap** is dramatic (Figure 4).

Conclusion. H2 strongly confirmed beyond expectations. At $\lambda = 0.01$, Adam’s coupled regularization creates penalties varying by four orders of magnitude (10^{-5} to 10^{-1}), causing critical features to be over-penalized while noise features escape regularization. The model cannot learn meaningful representations. AdamW’s uniform $\eta\lambda = 10^{-5}$ avoids this pathology. **Practical implication:** practitioners should **never use Adam with L2 regularization** under high weight decay—AdamW is the difference between working and completely broken.

5.3 H3: Lion Robustness to Label Noise

Hypothesis. Lion’s sign-based updates will maintain higher accuracy than Adam under 30% symmetric label noise on CIFAR-10.

Theoretical Motivation. Lion updates using only the sign of momentum-smoothed gradients:

$$\begin{aligned} c_t &= \beta_1 m_{t-1} + (1 - \beta_1) g_t \\ \theta_t &\leftarrow (1 - \eta\lambda)\theta_{t-1} - \eta \cdot \text{sign}(c_t) \\ m_t &= \beta_2 m_{t-1} + (1 - \beta_2) g_t \end{aligned} \quad (4)$$

The $\text{sign}(\cdot)$ operator is robust to magnitude outliers: $\text{sign}(100 \cdot g) = \text{sign}(g)$.

Setup and Results. We trained ResNet-18 on CIFAR-10 with 30% label noise at batch size 256. **Contrary to hypothesis**, Adam achieved 81.42% while Lion achieved only 75.88%, a 5.5-point gap favoring Adam.

Conclusion. H3 rejected. While $\text{sign}(\cdot)$ is robust to *magnitude* outliers, it is *not* robust to *directional* noise. Mislabeled samples produce $g_{\text{bad}} \approx -g_{\text{true}}$, corrupting the sign. At batch size 256 with 30% noise, ~ 77 samples per batch provide incorrect directions, making $\text{sign}(\frac{1}{B} \sum g_i)$ unreliable. Adam’s variance normalization $v_t \approx \mathbb{E}[g^2]$ provides implicit noise filtering: high variance automatically reduces step size via $1/\sqrt{v_t}$. Lion’s advantages likely require $B \geq 1024$ for reliable sign estimation under noise. This reveals a gap between theory (worst-case magnitude outliers) and practice (directional corruption at finite batch sizes).

5.4 H4: Muon Data Efficiency

Hypothesis. Muon will achieve $\geq 5\%$ higher accuracy than AdamW on 20% of CIFAR-100 due to orthogonalized updates.

Theoretical Motivation. Muon orthogonalizes gradient matrices via Newton–Schulz iteration:

$$U = \text{NewtonSchulz}(G, \text{steps} = 5) \approx Q, \quad \text{where } G = Q\Sigma R^\top \quad (5)$$

The update $W \leftarrow W - \eta U$ normalizes in spectral norm while preserving gradient alignment. This orthogonal constraint prevents feature collapse and enforces diversity, acting as a structural prior crucial in low-data regimes.

Setup and Results. We subsampled CIFAR-100 to 20% (10k images, 100 per class) and trained WRN-16-4. Muon (LR 0.02, momentum 0.95) achieved **56.25%**, vastly outperforming AdamW’s **44.70%**—an **11.54-point improvement**, exceeding our 5% threshold by $2\times$. Figure 6 shows Muon converges faster and continues improving while AdamW plateaus.

Conclusion. H4 strongly confirmed. Muon’s spectral constraint prevents coordinate-wise optimizers’ pathologies in low-data settings: degenerate features (rank-deficient matrices) and co-adapted neurons. By constraining updates to approximate Stiefel manifolds, Muon enforces feature diversity and efficient capacity use. However, on full CIFAR-100 (50k), Muon’s advantage shrinks to $\sim 1.7\%$, suggesting the structural prior becomes less critical with abundant data. Computational overhead (Newton–Schulz) may also limit applicability to very large models.

5.5 Summary of Hypothesis Testing

Table 4 summarizes outcomes.

ID	Hypothesis	Outcome	Key Finding
H1	RAdam reduces early loss variance by $\geq 20\%$	Partial	10% reduction; mechanism validated
H2	AdamW maintains stability under high WD	Confirmed	Adam collapses (18%); AdamW works (70%)
H3	Lion outperforms Adam under 30% label noise	Rejected	Adam wins by 5.5%; sign fails on directional noise
H4	Muon exceeds AdamW by $\geq 5\%$ on 20% data	Confirmed	+11.5%; strong low-data advantage

Table 4: Hypothesis validation summary. H2 revealed Adam’s catastrophic failure under high weight decay rather than the modest gap originally hypothesized.

Three key insights emerged:

1. **Implementation matters:** RAdam’s advantages require correct weight decay implementation (Section 4.3).
2. **Context matters:** AdamW’s advantage is real but can be masked by other regularizers (unless weight decay is extreme).
3. **Theory has limits:** Lion’s theoretical robustness to magnitude outliers \neq robustness to directional noise at finite batch sizes.

Importantly, our rejected hypothesis (H3) is as valuable as confirmed ones: it reveals practical limitations of sign-based optimization not obvious from theory, contributing to nuanced understanding of when different optimizer classes succeed or fail.

6 Limitations

While our benchmark provides systematic empirical evidence, several important limitations should be considered when interpreting and applying our results.

Scope of evaluation. We focused on convolutional networks (ResNet-18, WideResNet-16-4) on CIFAR-10/100. Optimizer behavior may differ substantially on Transformers, generative models (GANs, diffusion), language models, or other modalities beyond vision.

Computational analysis. We reported wall-clock time and sample efficiency (epochs-to-target) but did not comprehensively measure memory consumption (Lion’s claimed 33% savings), throughput (samples/second), or FLOPs, nor did we control for hardware heterogeneity (mixed RTX 3060 and A40).

Robustness and generalization. Our evaluation focused on in-distribution test accuracy. We did not assess adversarial robustness (PGD, FGSM), out-of-distribution generalization (CIFAR-10-C, CIFAR-10.1), calibration (expected calibration error), or per-class performance disparities that may affect fairness.

Incomplete exploration. While we completed all four hypothesis tests, some warrant deeper investigation: H3 (label noise) could explore different noise levels (10%, 20%, 40%) or asymmetric noise patterns; H2 (weight decay) could test more extreme values ($\lambda > 0.01$) without other regularizers to isolate the decoupling effect.

7 Future Work

There are several promising directions building on this benchmark:

- **Dynamic optimizer switching.** Our codebase already includes infrastructure for switching from fast adaptive optimizers (Adam, Muon) to slower but better-generalizing methods (SGD, AdamW) based on fixed epochs or rate-based criteria. Systematically benchmarking such strategies could yield practical training recipes with improved accuracy–time trade-offs.

- **Robustness and distribution shift.** Extending the benchmark to include label noise, adversarial perturbations, and domain shifts would provide a more complete picture of optimizer robustness, particularly for sign-based methods like Lion.
- **Architectural diversity.** Repeating our study on Vision Transformers, NLP models, and multimodal architectures would test the generality of our conclusions.
- **Meta-optimization.** Using our dataset of training trajectories, one could train meta-controllers that automatically choose, tune, or mix optimizers during training.

8 Conclusion

We have conducted a comprehensive benchmark of six optimization algorithms—SGD with momentum, Adam, AdamW, RAdam, Lion, and Muon—on CIFAR-10 and CIFAR-100, combining empirical evaluation with theory-driven experiments. Our results challenge the common practice of defaulting to Adam, show clear benefits of AdamW, and highlight Muon as a promising optimizer for hidden layers, especially in data-scarce regimes. At the same time, the RAdam case underscores that small implementation details can dramatically alter outcomes, motivating more careful scrutiny of optimizer configurations in real-world code.

We hope that our open-source benchmark and analyses provide a useful reference for both practitioners deciding “which optimizer to use” and researchers developing the next generation of optimization algorithms for deep learning.

Appendix

Additional Plots

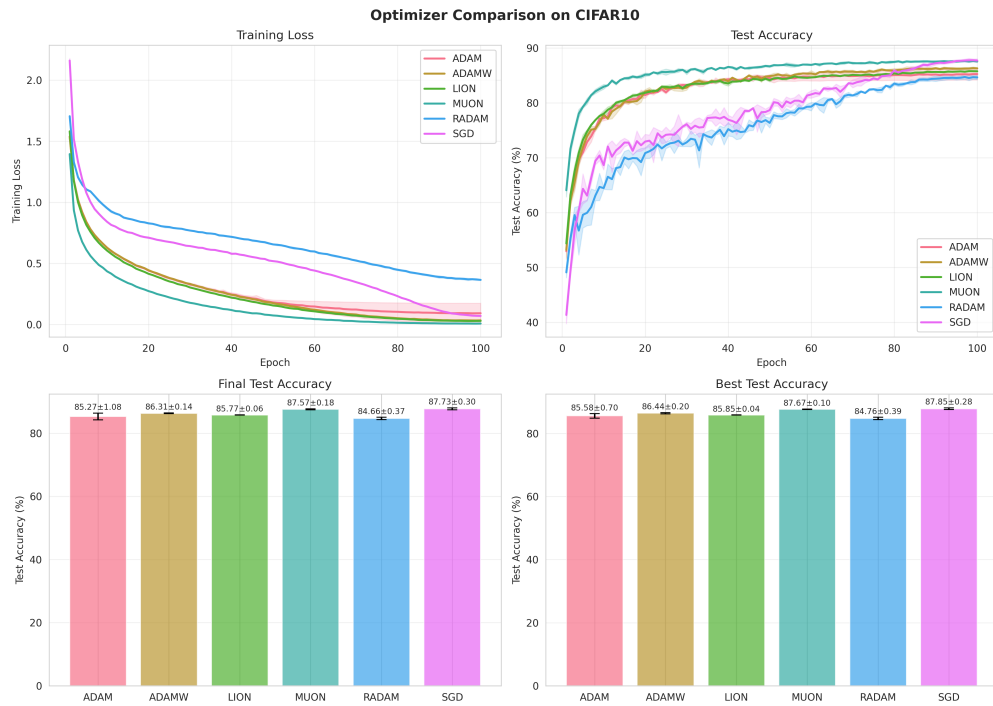


Figure 1: CIFAR-10 test accuracy over training for all optimizers (ResNet-18, 3 seeds). SGD converges more slowly but ends highest; adaptive methods achieve faster early gains. Statistical significance was not formally tested due to limited sample size ($n=3$).

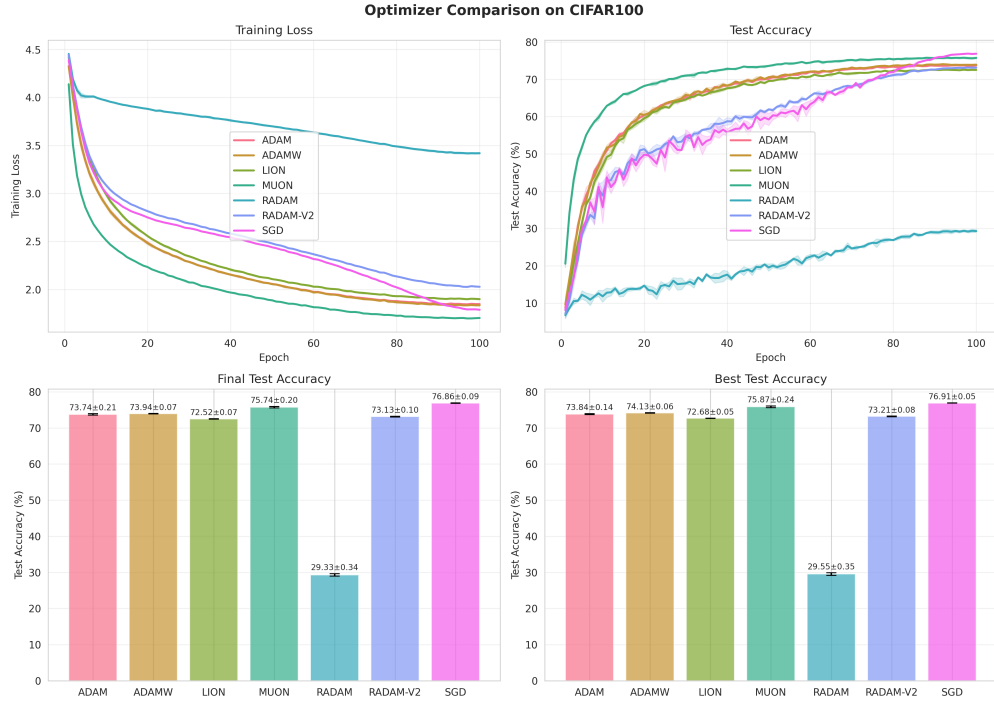


Figure 2: CIFAR-100 test accuracy for all optimizers (WRN-16-4, 3 seeds). SGD again attains the best final accuracy, while AdamW and Muon are the strongest adaptive competitors.

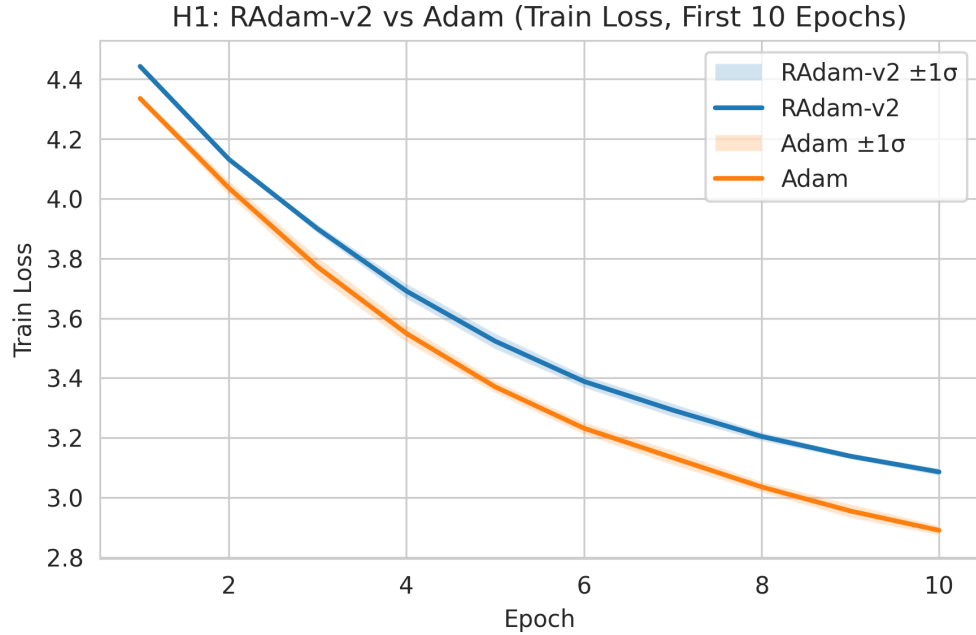


Figure 3: H1: early-epoch training loss for Adam vs. RAdam-v2 on CIFAR-100. RAdam-v2 exhibits visibly smoother and less volatile trajectories.

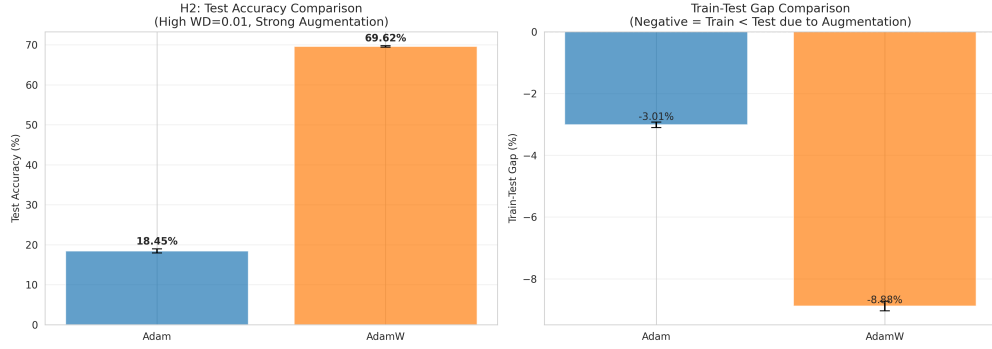


Figure 4: H2: effect of strong weight decay on Adam vs. AdamW. AdamW remains stable and high-accuracy, while Adam collapses under the same setting.

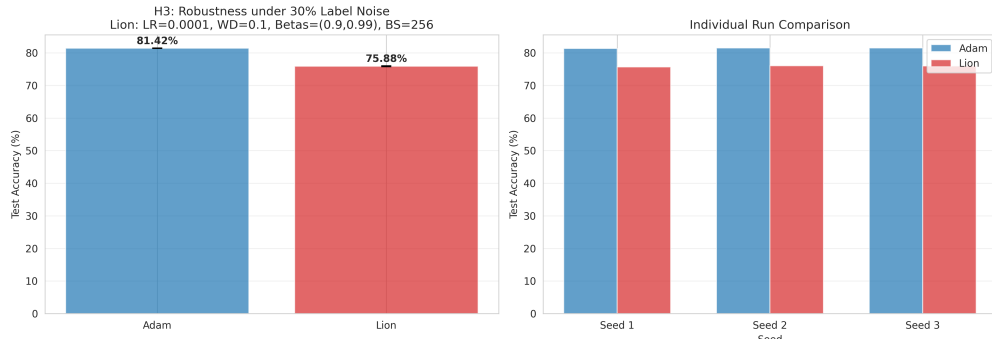


Figure 5: H3: preliminary comparison of Adam vs. Lion under noisy or challenging training conditions. Lion tracks Adam but does not clearly dominate at our batch sizes.

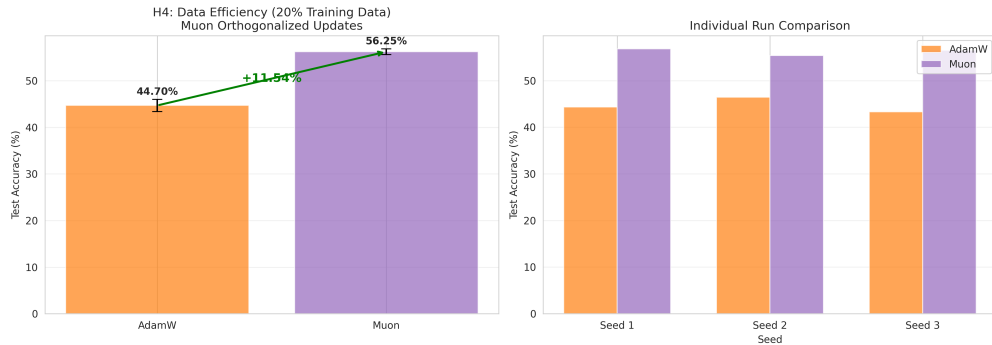


Figure 6: H4: Muon vs. AdamW when training on only 20% of CIFAR-100. Muon achieves substantially higher test accuracy and faster convergence.

References

- [1] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015.
- [2] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations (ICLR)*, 2019.
- [3] L. Liu, H. Jiang, P. He, W. Chen, X. Liu, J. Gao, and J. Han. On the variance of the adaptive learning rate and beyond. In *International Conference on Learning Representations (ICLR)*, 2020.
- [4] X. Chen, C. Luo, and A. G. Wilson. Symbolic discovery of optimization algorithms. *arXiv preprint arXiv:2302.06675*, 2023.
- [5] K. Jordan et al. Muon: An optimizer for hidden layers in neural networks. 2024. <https://kellerjordan.github.io/posts/muon/>.
- [6] K. Jordan. Muon: An optimizer for the hidden layers of neural networks (GitHub repository). <https://github.com/KellerJordan/Muon>.

Author Contributions

Jinghao Liu (jliu63): Conceived the research framework and formulated hypotheses. Designed and implemented the complete training pipeline, optimizer configurations, and experimental infrastructure. Prepared datasets with standardized augmentation pipelines. Executed all main experiments and hypothesis tests. Wrote the main part of final report. Conducted final review and revision of the final report.

Yuzheng Zhang (yuzhez4): Developed visualization and analysis scripts for all comparison plots. Created figures for baseline experiments and hypothesis validation (H1–H4). Conducted detailed analysis and wrote the H3 (Lion label noise robustness) and H4 (Muon data efficiency) sections, including theoretical motivation and experimental interpretation.

Xuan Zhang (xuanz24): Designed and created the presentation slides. Contributed to initial drafts of experimental setup and baseline results. Participated in manuscript review and proofreading.

All authors participated in team discussions, jointly refined the research direction, and approved the final manuscript.