

TRASMUON: TRUST-REGION ADAPTIVE SCALING FOR ORTHOGONALIZED MOMENTUM OPTIMIZERS

Anonymous authors

Paper under double-blind review

ABSTRACT

Muon-style optimizers leverage Newton-Schulz (NS) iterations to orthogonalize updates, yielding update geometries that often outperform Adam-series methods. However, this orthogonalization discards magnitude information, rendering training sensitive to step-size hyperparameters and vulnerable to high-energy bursts. To mitigate this, we introduce TrasMuon (Trust Region Adaptive Scaling Muon). TrasMuon preserves the near-isometric geometry of Muon while stabilizing magnitudes through (i) global RMS calibration and (ii) energy-based trust-region clipping. We demonstrate that while reintroducing adaptive scaling improves optimization efficiency, it typically exacerbates instability due to high-energy outliers. TrasMuon addresses this by defining a trust region based on relative energy ratios, confining updates to a stable zone. Empirical experiments on vision and language models demonstrate that TrasMuon converges faster than baselines. Furthermore, experiments without warmup stages confirm TrasMuon’s superior stability and robustness.

1 INTRODUCTION

Optimizer choice remains a bottleneck for training modern foundation models, shaping convergence and stability at scale (DeepSeek-AI et al., 2025; OpenAI, 2026; Anthropic, 2026). In practice, heterogeneous and heavy-tailed/outlier updates can trigger loss spikes and narrow the stable learning-rate region (Behrouz et al., 2025; Kimi Team et al., 2025; Park et al., 2025). Diagonal adaptive methods (Adam/AdamW and variants) provide robust coordinate-wise magnitude control (Kingma & Ba, 2017; Loshchilov & Hutter, 2019; Pagliardini et al., 2025; Marfinez, 2025), but do not exploit matrix-level update structure. Muon-style optimizers revisit matrix-structured updates via momentum orthogonalization (Bernstein, 2025; Jordan et al., 2024). However, orthogonalization mainly controls *geometry* and discards magnitude information, increasing sensitivity to step-size/warmup choices and vulnerability to bursty, axis-localized energy spikes (Behrouz et al., 2025; Kimi Team et al., 2025; Park et al., 2025).

We propose **TrasMuon** (Trust-Region Adaptive Scaling for **Muon**), which preserves Muon-style structured mixing while stabilizing magnitudes via global RMS calibration and feature-wise relative-energy trust-region damping. For a matrix parameter $W \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$, TrasMuon applies

$$\Delta W_t = -\hat{\eta}_t O_t^{\text{base}} \text{diag}(c_t), \quad c_t \in [c_{\min}, 1]^{d_{\text{in}}}. \quad (1)$$

Here O_t^{base} is obtained by NS orthogonalization plus lightweight row-wise second-moment scaling (NorMuon-style) (Li et al., 2025). The RMS-calibrated step size $\hat{\eta}_t$ improves cross-layer comparability and reduces step-size sensitivity (Bernstein & Newhouse, 2025; Large et al., 2024). The clipping vector c_t selectively suppresses high-energy feature axes while largely preserving the structured mixing factor, and we stabilize this signal via effective-time (schedule-free) weighting (Defazio et al., 2024).

Contributions: We introduce TRASMUON, which combines Muon-style mixing with feature-wise trust-region clipping. TrasMuon achieves faster early-stage convergence and improved stability, exhibiting reduced sensitivity to the learning-rate magnitude and scheduling choices. In particular, it consistently mitigates loss spikes induced by heavy-tailed, axis-localized gradient bursts, while preserving the characteristic geometry of the Muon-style optimizers.

2 RELATED WORK

Diagonal preconditioning and Adam-style optimizers. Adam/AdamW and many refinements remain default baselines due to robust diagonal second-moment scaling (Kingma & Ba, 2017; Loshchilov & Hutter, 2019; Yuan et al., 2025; Pagliardini et al., 2025; Marfinez, 2025; Shao et al., 2025; Gupta & Wojtowysch, 2025). These methods stabilize training via coordinate-wise magnitude control but do not explicitly exploit matrix-level structure; TrasMuon instead keeps matrix-structured mixing and adds trust-region adaptive scaling.

Beyond diagonal: block/matrix preconditioners and trust ratios. Richer preconditioners (e.g., K-FAC, Shampoo, Adafactor) capture non-diagonal curvature structure with tractable approximations (Martens & Grosse, 2015; Gupta et al., 2018; Shazeer & Stern, 2018). Layerwise norm/trust-ratio scaling such as LARS/LAMB controls step magnitudes by comparing parameter and update norms (You et al., 2017; 2020). TrasMuon does not estimate curvature factors; it constructs a near-isometric mixing factor via orthogonalization and stabilizes magnitudes using RMS calibration plus a trust-region clipping.

Orthogonalized directions and Muon-style updates. Muon-style optimizers use Newton–Schulz iterations to approximate polar factors, yielding near-isometric directions that can improve Transformer training (Jordan et al., 2024; Bernstein, 2025). Related perspectives connect orthogonalized updates to modular/geometry-aware optimization and practical variants (Bernstein & Newhouse, 2025; Large et al., 2024; Pethick et al., 2025; Ahn et al., 2025; Kumar et al., 2025; Khaled et al., 2025; Riabinin et al., 2025; Li et al., 2025). TrasMuon builds on these directions but targets a distinct failure mode: bursty, axis-localized energy spikes, addressed via feature-wise relative-energy clipping and temporal smoothing.

3 METHODOLOGY

We propose **TrasMuon (Trust-Region Adaptive Scaling for Muon)**, which *factorizes* matrix updates into a structured mixing factor (illustrated in Appendix B) and feature-wise trust-region adaptive scaling (in Algorithm 1).

Trust-region adaptive scaling. We control bursty, axis-localized updates using the per-column relative-energy ratio $r_{t,j} = E_{t,j} / (E_t^{\text{ref}} + \epsilon)$ and apply multiplicative damping, corresponding to an implicit trust-region constraint $r_{t,j} \lesssim \tau$ (tuned by α and optional trigger k). We measure column energy on pre-orthogonalization momentum M_t ,

$$E_{t,j} = \sum_{i=1}^{d_{\text{out}}} M_{t,ij}^2, \quad (2)$$

use a robust reference based on the median,

$$E_t^{\text{cur}} = \text{Quantile}_{0.5}(\{E_{t,j}\}), \quad (3)$$

$$E_t^{\text{ref}} = \beta_E E_{t-1}^{\text{ref}} + (1 - \beta_E) E_t^{\text{cur}}, \quad (4)$$

which resists inflation by sparse bursts (Hampel et al., 1986; Huber, 1981). We then apply a smooth, damping-only clip

$$c_{t,j}^{\text{raw}} = \frac{1}{1 + \alpha \log(1 + r_{t,j})}, \quad (5)$$

$$c_{t,j}^{\text{clip}} = \text{clip}(c_{t,j}^{\text{raw}}, c_{\min}, 1), \quad (6)$$

(optionally only when $r_{t,j} > k$), and set c_t by temporal smoothing (Appendix B).

Schedule-free temporal smoothing. We smooth the instantaneous clip with EMA

$$c_t^{\text{ema}} = \beta_c c_{t-1}^{\text{ema}} + (1 - \beta_c) c_t^{\text{inst}}, \quad (7)$$

and optionally apply schedule-free averaging (Defazio et al., 2024) (default $\gamma_t = \eta$):

$$S_t = S_{t-1} + \gamma_t^2, \quad C_t = C_{t-1} + \gamma_t^2 c_{t-1}^{\text{last}}, \quad (8)$$

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

$$c_t^{\text{avg}} = \frac{C_t}{S_t + \epsilon}, \quad c_t = (1 - \rho)c_t^{\text{ema}} + \rho c_t^{\text{avg}}. \tag{9}$$

We cache $c_t^{\text{last}} \leftarrow c_t$ between gate updates to avoid bias when the raw clip is computed every K steps. Convergence analysis can be seen in Appendix C.

4 EXPERIMENTS

TRASMUON is evaluated across four complementary settings: language-model pretraining under a fixed-budget protocol that highlights early- and late-stage dynamics and includes a minimal replication check, with full hyperparameters and additional plots reported in Appendix D; vision transformers trained with a standard ViT recipe on ImageNet-100 under multi-seed evaluation, together with an additional column-burst stress test, detailed in Appendix E; physics-informed neural networks for Helmholtz under controlled random-ROI sampling shifts that induce reproducible nonstationarity, alongside step-size alignment and learning-rate sensitivity diagnostics in Appendix F; and a controlled diagnostics study with column-localized outlier injection designed to probe the energy-indexed, feature-wise clipping mechanism and its boundary conditions, documented in Appendix G.

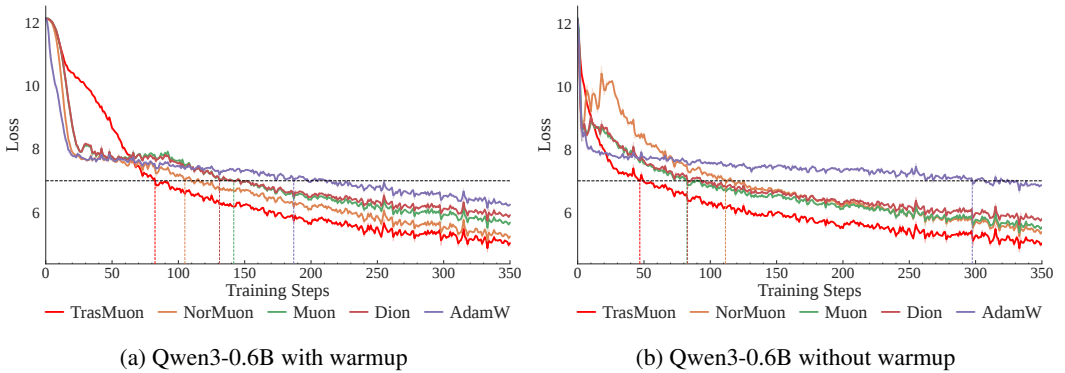


Figure 1: Early-stage training dynamics for Qwen3-0.6B from scratch (steps 0–350) under a warmup-stable-decay schedule, comparing (a) warmup-enabled runs and (b) warmup-free runs. The curves are smoothed using a time-weighted exponential moving average (EMA) with smoothing factor 0.1 for better visualization.

Training with warmup. Under the warmup-enabled schedule, all evaluated optimizers exhibit stable training in this fixed-hyperparameter pilot. Although TRASMUON is not the fastest within the first few dozen steps, it begins to reduce loss more rapidly after approximately 80 steps (illustrated in Fig. 1a). At a reference threshold of training loss (7.0), TRASMUON reaches the target in 80 steps, compared to 188 steps for AdamW (2.35×) and 140 steps for Muon (1.75×).¹

Training without warmup. Without warmup, optimizer becomes more sensitive to step-size calibration. Under the same shared learning rate and batch/sequence configuration, TRASMUON maintains a smooth loss trajectory in the early stage, while several baselines show larger loss oscillations as shown in Fig. 1b. At the same reference threshold of training loss, TRASMUON reaches the target in 48 steps, versus 298 for AdamW (6.21×) and 83 for Muon (1.73×). Despite the widespread use of warmup in pretraining, this setting remains meaningful, as the choice of warmup length is typically determined heuristically.

Late-stage behavior under a fixed budget. Beyond early-stage speed, TRASMUON also attains the lowest (or comparable-lowest) training loss in the late-stage window under both warmup settings in this pilot run. Late-stage curves and loss comparisons are reported in Fig. 2 of Appendix D.2. The early-stage gap narrows as training progresses into a slower-loss regime, suggesting that the benefit is strongest when training dynamics are most nonstationary.

¹We use loss = 7.0 as a representative checkpoint.

162 **Decreasing feature-wise energy concentration over training.** A plausible explanation is that
 163 early training exhibits stronger feature-wise anisotropy, where a small subset of hidden dimen-
 164 sions (columns) contributes disproportionately to gradient or momentum energy. In this regime,
 165 TRASMUON’s energy-based feature-wise clipping is engaged, selectively damping bursty columns
 166 and improving stability without discarding the structured Muon direction. As representations be-
 167 come better calibrated, energy may become more uniformly distributed across feature axes, reducing
 168 the prevalence of strongly concentrated column-localized bursts; consequently, the clipping signal
 169 weakens and the effective update becomes closer to the NorMuon backbone. This explanation is of-
 170 fered as a hypothesis and a direct characterization of activation/gradient anisotropy (e.g., via tracked
 171 energy ratios and gate statistics over time) is left to future work.

172 5 DISCUSSION AND LIMITATIONS

173 **What TRASMUON changes.** TRASMUON factorizes matrix updates into (i) a Muon-style near-
 174 isometric mixing factor constructed by Newton–Schulz orthogonalization and (ii) explicit magnitude
 175 controls: a global RMS-calibrated step size and a bounded, damping-only feature-wise clipping
 176 $c_{t,j} \in [c_{\min}, 1]$. This design targets a common practical tension: structured mixing can improve
 177 optimization geometry, while stable magnitudes govern learning-rate sensitivity and robustness to
 178 heavy-tailed bursts. When feature axes are semantically meaningful and bursts are axis-localized,
 179 TRASMUON selectively attenuates high-energy columns while largely preserving the Muon-style
 180 mixing structure.

181 **When feature-wise clipping and effective-time smoothing help.** Feature-wise clipping is most
 182 beneficial when update energy concentrates on a small subset of feature axes, reflected by large
 183 relative ratios $r_{t,j} = E_{t,j}/(E_t^{\text{ref}} + \epsilon)$. In this regime, multiplicative clipping suppresses burst-
 184 dominated columns without amplification. When clipping is recomputed sparsely (every K steps) or
 185 schedules vary, effective-time (schedule-free) averaging provides a stable long-horizon estimate by
 186 accumulating γ_t^2 -weighted statistics, reducing sensitivity to recomputation frequency and schedule
 187 details.

188 **Limitations.** (i) The formulation is most natural for 2D weight matrices; extending energy diag-
 189 nostics and damping to embeddings and higher-order tensors requires careful axis conventions. (ii)
 190 Newton–Schulz orthogonalization is sensitive to numerical precision; large-scale deployment bene-
 191 fits from precision-aware implementations. (iii) TRASMUON introduces additional hyperparameters
 192 and design choices (e.g., $K, \alpha, k, c_{\min}, \rho$ and the robust reference), and their interactions with model
 193 scale and data regimes merit broader sweeps.

194 6 CONCLUSION

195 We presented TRASMUON, a Muon-style optimizer that combines (i) NS-based near-isometric
 196 mixing factors with (ii) explicit magnitude stabilization via global RMS calibration and bounded,
 197 feature-wise trust-region clipping, optionally smoothed by effective-time (schedule-free) averaging.
 198 Across the evaluated workloads, TRASMUON improves training stability and achieves competitive
 199 or better final performance than strong baselines. Controlled diagnostics further support the intended
 200 mechanism: column-localized energy bursts increase relative energy ratios and lead to stronger ap-
 201 plied damping, while ablations (e.g., NOCLIP) and broken-axis settings help rule out trivial expla-
 202 nations such as uniform step-size reduction. On practical tasks including language-model training,
 203 vision transformers, and PINNs under ROI-induced sampling shifts, TRASMUON yields faster or
 204 more stable optimization dynamics and improved robustness (Future Work can be deferred in Ap-
 205 ppendix A).

206 IMPACT STATEMENT

207 This paper presents work whose goal is to advance the field of machine learning. There are many
 208 potential societal consequences of our work, none of which we feel must be specifically highlighted
 209 here.

REFERENCES

- 216
217
218 Kwangjun Ahn, Byron Xu, Natalie Abreu, Ying Fan, Gagik Magakyan, Pratyusha Sharma,
219 Zheng Zhan, and John Langford. Dion: Distributed orthonormalized updates. *arXiv preprint*
220 *arXiv:2504.05295*, 2025.
- 221 Ambityga. ImageNet100. [https://www.kaggle.com/datasets/ambityga/](https://www.kaggle.com/datasets/ambityga/imagenet100)
222 [imagenet100](https://www.kaggle.com/datasets/ambityga/imagenet100), 2021. Accessed: 2026-01-23.
- 223
224 Anthropic. Introducing Claude Opus 4.5, 2026. URL [https://www.anthropic.com/news/](https://www.anthropic.com/news/claude-opus-4-5)
225 [claude-opus-4-5](https://www.anthropic.com/news/claude-opus-4-5). Accessed: 2026-01-12.
- 226 Ali Behrouz, Zeman Li, Praneeth Kacham, Majid Daliri, Yuan Deng, Peilin Zhong, Meisam Raza-
227 viyayn, and Vahab Mirrokni. ATLAS: Learning to optimally memorize the context at test time.
228 *arXiv preprint arXiv:2505.23735*, 2025.
- 229
230 Jeremy Bernstein. Deriving Muon, 2025. URL [https://jeremybernste.in/writing/](https://jeremybernste.in/writing/deriving-muon)
231 [deriving-muon](https://jeremybernste.in/writing/deriving-muon).
- 232
233 Jeremy Bernstein and Laker Newhouse. Modular duality in deep learning. In *International Confer-*
234 *ence on Machine Learning*, 2025.
- 235
236 DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, et al. DeepSeek-V3 technical report.
arXiv preprint arXiv:2412.19437, 2025.
- 237
238 Aaron Defazio, Xingyu Alice Yang, Ahmed Khaled, Konstantin Mishchenko, Harsh Mehta, and
239 Ashok Cutkosky. The road less scheduled. In *Advances in Neural Information Processing Sys-*
240 *tems*, 2024.
- 241
242 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
243 hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition*,
2009.
- 244
245 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
246 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszko-
247 reit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at
scale. In *International Conference on Learning Representations*, 2021.
- 248
249 Zhiwei Gao, Liang Yan, and Tao Zhou. Failure-informed adaptive sampling for PINNs. *SIAM*
250 *Journal on Scientific Computing*, 45(4):A1971–A1994, 2023.
- 251
252 Kanan Gupta and Stephan Wojtowytsch. Nesterov acceleration in benignly non-convex landscapes.
253 In *The International Conference on Learning Representations*, 2025.
- 254
255 Vineet Gupta, Tomer Koren, and Yoram Singer. Shampoo: Preconditioned stochastic tensor opti-
mization. In *International Conference on Machine Learning*, 2018.
- 256
257 Frank R. Hampel, Elvezio M. Ronchetti, Peter J. Rousseeuw, and Werner A. Stahel. *Robust Statis-*
258 *tics: The Approach Based on Influence Functions*. John Wiley & Sons, 1986.
- 259
Peter J. Huber. *Robust Statistics*. John Wiley & Sons, 1981.
- 260
261 Keller Jordan, Yuchen Jin, Vlado Boza, Jiacheng You, Franz Cesista, Laker Newhouse, and Jeremy
262 Bernstein. Muon: An optimizer for hidden layers in neural networks, 2024. URL [https:](https://kellerjordan.github.io/posts/muon/)
263 [//kellerjordan.github.io/posts/muon/](https://kellerjordan.github.io/posts/muon/).
- 264
265 Ahmed Khaled, Kaan Ozkara, Tao Yu, Mingyi Hong, and Youngsuk Park. MuonBP: Faster Muon
via block-periodic orthogonalization. *arXiv preprint arXiv:2510.16981*, 2025.
- 266
267 Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, et al.
268 Kimi K2: Open agentic intelligence. *arXiv preprint arXiv:2507.20534*, 2025.
- 269
Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
arXiv:1412.6980, 2017.

- 270 Bhavesh Kumar, Roger Jin, and Jeffrey Quesnelle. CurvaDion: Curvature-adaptive distributed or-
271 thonormalization. *arXiv preprint arXiv:2512.13728*, 2025.
- 272
- 273 Tim Large, Yang Liu, Minyoung Huh, Hyojin Bahng, Phillip Isola, and Jeremy Bernstein. Scalable
274 optimization in the modular norm. In *Advances in Neural Information Processing Systems*, 2024.
- 275
- 276 Zichong Li, Liming Liu, Chen Liang, Weizhu Chen, and Tuo Zhao. NorMuon: Making Muon more
277 efficient and scalable. *arXiv preprint arXiv:2510.05491*, 2025.
- 278
- 279 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-
280 ence on Learning Representations*, 2019.
- 281
- 282 Anton Lozhkov, Loubna Ben Allal, Leandro von Werra, and Thomas Wolf. FineWeb-Edu: The finest
283 collection of educational content, 2024. URL [https://huggingface.co/datasets/
284 HuggingFaceFW/fineweb-edu](https://huggingface.co/datasets/HuggingFaceFW/fineweb-edu). Accessed: 2026-01-01.
- 285
- 286 Mitchell Marfinez. Evolving Deep Learning Optimizers. *arXiv preprint arXiv:2512.11853*, 2025.
- 287
- 288 James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate
289 curvature. In *International Conference on Machine Learning*, 2015.
- 290
- 291 OpenAI. Introducing GPT-5.2, 2026. URL [https://openai.com/index/
292 introducing-gpt-5-2/](https://openai.com/index/introducing-gpt-5-2/). Accessed: 2026-01-12.
- 293
- 294 Matteo Pagliardini, Pierre Ablin, and David Grangier. The AdEMAMix optimizer: Better, faster,
295 older. In *International Conference on Learning Representations*, 2025.
- 296
- 297 Jungwoo Park, Taewhoo Lee, Chanwoong Yoon, Hyeon Hwang, and Jaewoo Kang. Outlier-safe pre-
298 training for robust 4-Bit quantization of large language models. *arXiv preprint arXiv:2506.19697*,
299 2025.
- 300
- 301 Thomas Pethick, Wanyun Xie, Kimon Antonakopoulos, Zhenyu Zhu, Antonio Silveti-Falls, and
302 Volkan Cevher. Training deep learning models with norm-constrained LMOs. In *International
303 Conference on Machine Learning*, 2025.
- 304
- 305 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
306 models are unsupervised multitask learners. *OpenAI blog*, 2019.
- 307
- 308 Artem Riabinin, Egor Shulgin, Kaja Grutkowska, and Peter Richtárik. Gluon: Making Muon &
309 Scion great again! (bridging theory and practice of LMO-based optimizers for LLMs). *arXiv
310 preprint arXiv:2505.13416*, 2025.
- 311
- 312 Yichuan Shao, Shiqian Weng, Haijing Sun, Qian Gao, Le Zhang, Zhiqiang Mao, Shuai Xu, Zhitao
313 Zhang, and Lei Xing. BDS-Adam optimizer integrating adaptive variance rectification with semi-
314 adaptive gradient smoothing. *Scientific Reports*, 15(1):36906, 2025.
- 315
- 316 Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost.
317 In *International Conference on Machine Learning*, 2018.
- 318
- 319 Shashank Subramanian, Robert M. Kirby, Michael W. Mahoney, and Amir Gholami. Adaptive self-
320 supervision algorithms for physics-informed neural networks. *arXiv preprint arXiv:2207.04084*,
321 2022.
- 322
- 323 Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and
Hervé Jégou. Training data-efficient image transformers & distillation through attention. In
International Conference on Machine Learning, 2021.
- 324
- 325 Chenxi Wu, Min Zhu, Qinyang Tan, Yadhu Kartha, and Lu Lu. A comprehensive study of non-
326 adaptive and residual-based adaptive sampling for physics-informed neural networks. *Computer
327 Methods in Applied Mechanics and Engineering*, 403:115671, 2023.
- 328
- 329 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,
330 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint
331 arXiv:2505.09388*, 2025.

324 Yang You, Igor Gitman, and Boris Ginsburg. Large batch training of convolutional networks. *arXiv*
325 *preprint arXiv:1708.03888*, 2017.
326

327 Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan
328 Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep
329 learning: Training bert in 76 minutes. In *International Conference on Learning Representations*,
330 2020.

331 Huizhuo Yuan, Yifeng Liu, Shuang Wu, Xun Zhou, and Quanquan Gu. MARS: Unleashing the
332 power of variance reduction for training large models. In *International Conference on Machine*
333 *Learning*, 2025.
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

A FUTURE WORK

- **Theory and guarantees.** Develop stability and convergence analyses for orthogonalized updates combined with bounded, damping-only feature-wise clipping. A promising direction is to formalize a Lyapunov-style descent argument under explicit alignment conditions and bounded-update properties, and to connect empirical spectral diagnostics (e.g., the effective update statistic $A_{\text{eff}} = W^\top \Delta W$) to provable stability regimes.
- **Generalizing feature axes beyond 2D matrices.** Extend EnergyCol-style clipping to convolutional kernels, embeddings, and higher-order tensors by defining principled “feature axes” (e.g., input-channel, attention-head, or group dimensions). We also plan to explore block-wise and low-rank variants that preserve interpretability while reducing per-step overhead.
- **Adaptive burst modeling with transparent control.** Replace fixed clipping hyperparameters (e.g., c_{\min} , α , update period K) with lightweight, interpretable adaptations driven by online tail statistics of the energy distribution (quantiles, kurtosis, robust outlier scores), while maintaining the damping-only constraint and avoiding hidden amplification.
- **Systems and numerical precision.** Improve the efficiency and robustness of Newton–Schulz orthogonalization in mixed precision and distributed settings. This includes precision-aware kernels, communication-efficient implementations, and amortized/approximate orthogonalization strategies that retain most of the directional benefit at lower cost.
- **Broader robustness regimes and downstream impact.** Evaluate TRASMUON under realistic forms of nonstationarity common in large-scale training (curriculum shifts, domain-mixture changes, sequence-length spikes, and data-quality transitions), and study how clipping statistics correlate with downstream robustness and generalization.

B TRASMUON ALGORITHM (EXTENDED)

TrasMuon explicitly decouples *update geometry* (direction) from *step-size control* (magnitude). For a matrix parameter $W \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ with stochastic gradient $G_t = \nabla_W \mathcal{L}(W_t)$, it applies multiplicative coupling as

$$\Delta W_t = -\hat{\eta}_t O_t^{\text{base}} \text{diag}(c_t), \quad (10)$$

where $O_t^{\text{base}} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ provides a structured direction, $\hat{\eta}_t$ is a row-wise RMS-calibrated global step size, and $c_t \in [c_{\min}, 1]^{d_{\text{in}}}$ is a *feature-axis* (column-wise) damping vector (damping-only, no amplification).

B.1 ORTHOGONALIZED DIRECTIONS VIA NEWTON–SCHULZ

TrasMuon maintains an exponential moving average of gradients

$$M_t = \beta_1 M_{t-1} + (1 - \beta_1) G_t. \quad (11)$$

To extract a rotation-robust, near-isometric direction, we approximate the polar factor of M_t . For numerical stability of Newton–Schulz (NS) iterations, we remove the scale gauge by RMS-normalizing

$$\tilde{M}_t = \frac{M_t}{\|M_t\|_F / \sqrt{d_{\text{out}} d_{\text{in}}} + \epsilon}, \quad (12)$$

and apply T NS steps to obtain

$$O_t \approx \text{NS}(\tilde{M}_t; T) \approx \tilde{M}_t (\tilde{M}_t^\top \tilde{M}_t)^{-1/2}, \quad (13)$$

yielding a structured direction that is less sensitive to axis rotations than elementwise or diagonal preconditioning (e.g., Muon-style orthogonalized updates).

Algorithm 1 TRASMUON: Muon + Adaptive Scaling + Trust Region + Schedule-Free Smoothed

432
433
434 **Input:** $W \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$, base lr η , $\beta_1, \beta_2, \epsilon$, NS steps T , weight decay λ , $c_{\text{min}}, \alpha, \beta_E, \beta_c$, trigger k , update
435 period K , warmup T_w , mix ρ
436 Initialize $M \leftarrow 0$, $v^{\text{row}} \leftarrow 0$, $E^{\text{ref}} \leftarrow 0$, $c^{\text{ema}} \leftarrow 1$, $c^{\text{last}} \leftarrow 1$, $S \leftarrow 0$, $C \leftarrow 0$, $\gamma_t \leftarrow \eta$
437 **repeat**
438 $G \leftarrow \nabla_W \mathcal{L}(W)$; $W \leftarrow (1 - \eta\lambda)W$
439 **Momentum:** $M \leftarrow \beta_1 M + (1 - \beta_1)G$
440 **Orthogonalized direction:** $O \leftarrow \text{NS}(\tilde{M}; T)$
441 $v^{\text{row}} \leftarrow \beta_2 v^{\text{row}} + (1 - \beta_2) \text{mean}_j(O_{:,j}^2)$
442 $O^{\text{base}} \leftarrow \text{diag}((v^{\text{row}} + \epsilon)^{-1/2}) O$
443 **Calibration:** $\hat{\eta} \leftarrow \eta \frac{\sqrt{d_{\text{out}} d_{\text{in}}}}{\|O^{\text{base}}\|_F + \epsilon}$
444 **Column energy:** $E_j \leftarrow \sum_i M_{ij}^2$
445 **Robust reference:** $E^{\text{cur}} \leftarrow \text{Quantile}_{0.5}(\{E_j\})$
446 **EMA smooth:** $E^{\text{ref}} \leftarrow \beta_E E^{\text{ref}} + (1 - \beta_E) E^{\text{cur}}$
447 **Schedule-free accumulators:** $S \leftarrow S + \gamma_t^2$, $C \leftarrow C + \gamma_t^2 c^{\text{last}}$
448 $c^{\text{avg}} \leftarrow C / (S + \epsilon)$
449 **if** $t > T_w$ **and** $t \bmod K = 0$ **then**
450 $r_j \leftarrow E_j / (E^{\text{ref}} + \epsilon)$
451 $c_j^{\text{clip}} \leftarrow \text{clip}\left(\frac{1}{1 + \alpha \log(1 + r_j)}, c_{\text{min}}, 1\right)$
452 **Trigger (optional):** $c_j^{\text{inst}} \leftarrow c_j^{\text{clip}}$ if $r_j > k$ else 1
453 **EMA smooth:** $c^{\text{ema}} \leftarrow \beta_c c^{\text{ema}} + (1 - \beta_c) c^{\text{inst}}$
454 **end if**
455 **Long-Short term Mixing:** $c \leftarrow 1$ if $t \leq T_w$ else $(1 - \rho)c^{\text{ema}} + \rho c^{\text{avg}}$
456 $c^{\text{last}} \leftarrow c$ {cached between updates}
457 **Update:** $W \leftarrow W - \hat{\eta} (O^{\text{base}} \odot \text{ExpandCols}(c))$
458 **until** training ends

B.2 ROW-SECOND-MOMENT SCALING AND RMS-CALIBRATED STEP SIZE

461
462 Orthogonalization primarily shapes *direction*. To stabilize *magnitude* across layers and time, we
463 apply lightweight row-wise second-moment scaling (as in NorMuon (Li et al., 2025)):
464

$$465 \quad v_t^{\text{row}} = \beta_2 v_{t-1}^{\text{row}} + (1 - \beta_2) \text{mean}_j(O_{t,j}^{\odot 2}), \quad (14)$$

$$466 \quad O_t^{\text{base}} = \text{diag}((v_t^{\text{row}} + \epsilon)^{-1/2}) O_t. \quad (15)$$

467 Row scaling addresses row-wise heterogeneity, while a *global* row-wise calibration controls the
468 update norm. We set

$$470 \quad \hat{\eta}_t = \eta \cdot \frac{\sqrt{d_{\text{out}} d_{\text{in}}}}{\|O_t^{\text{base}}\|_F + \epsilon}, \quad (16)$$

471 so that the per-step RMS magnitude of ΔW_t is on the order of η . Since $c_t \leq 1$ elementwise, equa-
472 tion 16 also implies an explicit Frobenius-norm bound $\|\Delta W_t\|_F \leq \eta \sqrt{d_{\text{out}} d_{\text{in}}}$ (up to ϵ), reducing
473 sensitivity to layer shape and fluctuations in the orthogonalized direction (Bernstein & Newhouse,
474 2025; Large et al., 2024).

B.3 ENERGY-BASED FEATURE-WISE TRUST-REGION CLIPPING

478 **Motivation.** In practice, instability often arises from *bursty magnitudes* that concentrate on a small
479 subset of feature axes (columns), causing loss spikes and narrowing the stable learning-rate region.
480 TrasMuon therefore introduces *feature-wise clipping*: it dampens only the high-energy feature di-
481 rections while preserving the Muon-like direction structure in O_t^{base} .

482
483 **Column energy and a robust reference.** The column energy is measured from momentum M_t :

$$484 \quad E_{t,j} = \sum_{i=1}^{d_{\text{out}}} M_{t,ij}^2, \quad j = 1, \dots, d_{\text{in}}. \quad (17)$$

486 We summarize the typical energy level at step t by a quantile statistic

$$487 E_t^{\text{cur}} = \text{Quantile}_q(\{E_{t,j}\}_{j=1}^{d_{\text{in}}}), \quad q = 0.5, \quad (18)$$

489 and maintain a running reference via an EMA updated every step:

$$490 E_t^{\text{ref}} = \beta_E E_{t-1}^{\text{ref}} + (1 - \beta_E) E_t^{\text{cur}}. \quad (19)$$

491 Using a quantile (median) yields a high-breakdown reference: the sample median has a 50% break-
492 down point, so a sparse set of high-energy columns cannot arbitrarily inflate E_t^{ref} and thereby “move
493 the clipping threshold” in response to the outliers being clipped (Hampel et al., 1986; Huber, 1981).
494

495 **Relative ratio and clipping-style damping.** We define a dimensionless ratio as follow:

$$496 r_{t,j} = \frac{E_{t,j}}{E_t^{\text{ref}} + \epsilon}. \quad (20)$$

499 A hard energy cap $E_{t,j} \leq k E_t^{\text{ref}}$ gives the column-wise analogue of norm clipping:

$$500 c_{t,j}^{\text{hard}} = \min\left(1, \sqrt{\frac{k E_t^{\text{ref}}}{E_{t,j} + \epsilon}}\right), \quad (21)$$

503 which enforces $r_{t,j} \leq k$ after rescaling. Trasmuon leverages a smooth, numerically stable *soft*
504 *clipping* rule:

$$505 c_{t,j}^{\text{raw}} = \frac{1}{1 + \alpha \log(1 + r_{t,j})}, \quad c_{t,j}^{\text{gate}} = \text{clip}(c_{t,j}^{\text{raw}}, c_{\min}, 1). \quad (22)$$

507 This rule is bounded and avoids power-law instabilities as $r_{t,j} \rightarrow 0$. Importantly, $c_{t,j} \leq 1$ for all j ,
508 so the mechanism is strictly damping-only and can be interpreted as a trust-region safety mechanism
509 in *feature space*.
510

511 **Triggered vs. continuous clipping.** Optionally, damping can be also applied only when $r_{t,j}$ ex-
512 ceeds a triggering threshold k :

$$513 \bar{c}_{t,j} = \begin{cases} c_{t,j}^{\text{gate}}, & r_{t,j} > k, \\ 1, & \text{otherwise,} \end{cases} \quad (23)$$

516 so that non-burst columns remain unchanged and the gate acts as an event-driven clip.

518 B.4 TEMPORAL SMOOTHING AND SCHEDULE-FREE AVERAGING

519 To reduce short-term noise and avoid sensitivity to the gate-update period, we smooth the instanta-
520 neous clip by first applying EMA smoothing:

$$521 c_t^{\text{ema}} = \beta_c c_{t-1}^{\text{ema}} + (1 - \beta_c) c_t^{\text{inst}}, \quad (24)$$

523 where c_t^{inst} denote the clip applied at step t (either c_t^{gate} or \bar{c}_t depending on triggering). Second, we
524 maintain a schedule-free average using an effective step weight γ_t (default $\gamma_t = \eta$) (Defazio et al.,
525 2024). Define the scalar accumulator $S_t \in \mathbb{R}$ and vector accumulator $C_t \in \mathbb{R}^{d_{\text{in}}}$:

$$526 S_t = S_{t-1} + \gamma_t^2, \quad C_t = C_{t-1} + \gamma_t^2 c_{t-1}^{\text{last}}, \quad c_t^{\text{avg}} = \frac{C_t}{S_t + \epsilon}, \quad (25)$$

528 where c_{t-1}^{last} is the most recently applied clip (cached between updates). We then mix short- and
529 long-term estimates:

$$530 c_t = (1 - \rho) c_t^{\text{ema}} + \rho c_t^{\text{avg}}, \quad c_t^{\text{last}} \leftarrow c_t. \quad (26)$$

532 This effective-time averaging reduces sensitivity to warmup length and total training steps, and
533 prevents bias when the raw clip is computed only every K steps.

534 B.5 FINAL TRASMUON UPDATE

536 Substituting the components into equation 10 yields the final update

$$537 \Delta W_t = -\hat{\eta}_t O_t^{\text{base}} \text{diag}(c_t), \quad (27)$$

538 which preserves Muon-like directional geometry via O_t^{base} while controlling bursty magnitudes
539 along feature axes through damping-only feature clipping c_t .

C CONVERGENCE ANALYSIS

Scope. We provide a convergence *framework* for TrasMuon/energy clipping updates, separating (i) unconditional algebraic properties (bounded update norm; damping-only contraction) from (ii) mild alignment assumptions connecting the structured update to descent.

Damping-only contraction. For any matrix A and any $c \in [0, 1]^n$, right-multiplication by $\text{diag}(c)$ cannot increase the Frobenius norm:

$$\|A \text{diag}(c)\|_F \leq \|A\|_F. \quad (28)$$

RMS calibration. With $\hat{\eta}_t = \eta\sqrt{mn}/(\|O_t^{\text{base}}\|_F + \epsilon)$ and damping-only $c_t \leq \mathbf{1}$, the update norm is uniformly bounded:

$$\|\Delta W_t\|_F \leq \eta\sqrt{mn} \quad \forall t, \quad (29)$$

independently of transient gradient spikes.

Stationarity under smoothness and alignment. Under standard L -smoothness and a mild alignment condition (Appendix C), TrasMuon satisfies an expected first-order stationarity bound of the form

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(W_t)\|_F^2 \leq \frac{\mathbb{E}[f(W_0)] - f^*}{\mu \eta T} + \frac{L}{2\mu} \eta mn, \quad (30)$$

for a constant $\mu > 0$ capturing effective descent. Importantly, the EMA/schedule-free construction of c_t only affects how the clip is computed, while the theory relies solely on the invariant $c_{t,j} \in [c_{\min}, 1]$.

C.1 ALGORITHMIC ABSTRACTION

Let $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ and let G_t be a stochastic gradient at W_t . We study updates of the form

$$W_{t+1} = (1 - \eta\lambda) W_t + \Delta W_t, \quad \Delta W_t := -\hat{\eta}_t U_t, \quad U_t := O_t^{\text{base}} \text{diag}(c_t), \quad (31)$$

with RMS-calibrated step size

$$\hat{\eta}_t = \eta \cdot \frac{\sqrt{mn}}{\|O_t^{\text{base}}\|_F + \epsilon}, \quad (32)$$

and damping-only clipping

$$0 < c_{\min} \leq c_{t,j} \leq 1 \quad \forall j, t. \quad (33)$$

C.2 DETERMINISTIC ALGEBRAIC PROPERTIES

Lemma C.1 (Damping-only contraction). *For any $A \in \mathbb{R}^{m \times n}$ and any $c \in [0, 1]^n$, $\|A \text{diag}(c)\|_F \leq \|A\|_F$.*

Proof. $\|A \text{diag}(c)\|_F^2 = \sum_{j=1}^n c_j^2 \|A_{\cdot j}\|_2^2 \leq \sum_{j=1}^n \|A_{\cdot j}\|_2^2 = \|A\|_F^2. \quad \square$

Lemma C.2 (Row-wise RMS calibration). *Under equation 31–equation 33, for all t ,*

$$\|\Delta W_t\|_F \leq \eta\sqrt{mn}. \quad (34)$$

Proof. By Lemma C.1, $\|U_t\|_F = \|O_t^{\text{base}} \text{diag}(c_t)\|_F \leq \|O_t^{\text{base}}\|_F$. Thus

$$\|\Delta W_t\|_F = \hat{\eta}_t \|U_t\|_F \leq \eta\sqrt{mn} \cdot \frac{\|O_t^{\text{base}}\|_F}{\|O_t^{\text{base}}\|_F + \epsilon} \leq \eta\sqrt{mn}.$$

\square

594 C.3 ASSUMPTIONS FOR DESCENT

595 **Assumption C.3** (*L-smoothness*). f is L -smooth with respect to Frobenius norm: $\|\nabla f(X) - \nabla f(Y)\|_F \leq L\|X - Y\|_F$.

598 **Assumption C.4** (*Stochastic gradients*). $\mathbb{E}[G_t | W_t] = \nabla f(W_t)$ and $\mathbb{E}[\|G_t - \nabla f(W_t)\|_F^2 | W_t] \leq \sigma^2$.

600 **Assumption C.5** (*Alignment on the realized update*). There exists $\mu_\Delta > 0$ such that for all t ,

$$601 \mathbb{E}[\langle \nabla f(W_t), \Delta W_t \rangle | W_t] \leq -\mu_\Delta \eta \|\nabla f(W_t)\|_F^2, \quad (35)$$

602 where $\Delta W_t = -\hat{\eta}_t U_t$.

605 C.4 EXPECTED STATIONARITY

607 **Lemma C.6** (*Smoothness descent*). Under Assumption C.3, for any Δ ,

$$609 f(W_t + \Delta) \leq f(W_t) + \langle \nabla f(W_t), \Delta \rangle + \frac{L}{2} \|\Delta\|_F^2. \quad (36)$$

611 **Theorem C.7** (*Expected stationarity for RMS-calibrated, damping-only updates*). Assume C.3 and C.5 with $\lambda = 0$. Then for any $T \geq 1$,

$$614 \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(W_t)\|_F^2 \leq \frac{\mathbb{E}[f(W_0)] - f^*}{\mu_\Delta \eta T} + \frac{L}{2\mu_\Delta} \eta mn, \quad (37)$$

616 where $f^* = \inf_W f(W)$.

618 *Proof.* Apply Lemma C.6 with $\Delta = \Delta W_t$, take conditional expectation, use Assumption C.5, and bound $\|\Delta W_t\|_F^2$ by Lemma C.2. Sum over $t = 0, \dots, T - 1$ and telescope $f(W_t)$. \square

622 C.5 EXTENSIONS: WEIGHT DECAY, STOCHASTICITY, AND PL

623 **Weight decay.** Decoupled weight decay can be handled by analyzing $f_\lambda(W) = f(W) + \frac{\lambda}{2} \|W\|_F^2$ or treating $(1 - \eta\lambda)W_t$ as an additional contraction term.

626 **Stochastic gradients.** Under Assumption C.4, the bound acquires an additional $\mathcal{O}(\eta\sigma^2)$ term as in standard SGD analyses.

629 **PL / strong convexity.** If f satisfies the PL condition, one obtains linear-type convergence to an $\mathcal{O}(\eta mn)$ (or $\mathcal{O}(\eta\sigma^2)$) neighborhood.

632 D SUPPLEMENTARY RESULTS: LANGUAGE-MODEL PRETRAINING

634 **Scope.** This appendix provides supplementary visualizations for language-model pretraining and a minimal replication check. The intent is *not* to introduce new claims beyond the main text, but to (i) document late-stage behavior under the fixed training budget and (ii) reduce the risk that qualitative trends are specific to a single random seed.

639 D.1 EXPERIMENT SETTINGS

641 **Fixed-budget protocol and reporting (language model pretraining).** TRASMUON is evaluated in a controlled, short-run pretraining-style setting and compared against four baseline optimizers: AdamW (Loshchilov & Hutter, 2019), Muon (Jordan et al., 2024), Dion (Ahn et al., 2025), and NorMuon (Li et al., 2025). Decoder-only Transformer models are trained from random initialization, including GPT-2 (Radford et al., 2019) and Qwen3-0.6B (Yang et al., 2025), on FineWeb-Edu (Lozhkov et al., 2024). All runs follow the same fixed-budget protocol as Section 4: training proceeds for 1500 optimization steps with sequence length 1024 and a fixed global (effective) batch size of 1024, corresponding to $1500 \times 1024 \times 1024 \approx 1.57 \times 10^9$ training tokens. Unless

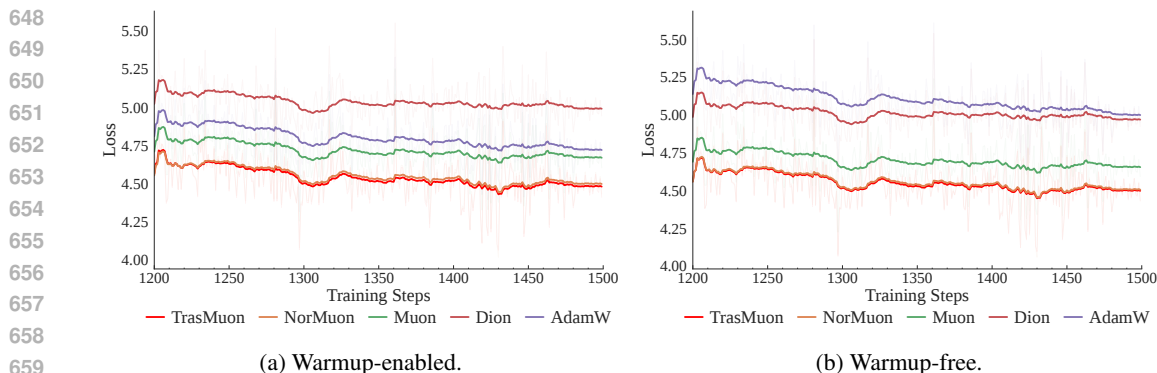


Figure 2: **Qwen3-0.6B late-stage window (steps 1200–1500) under a fixed-budget protocol.** Training loss under warmup-enabled vs. warmup-free variants of the same warmup–stable–decay schedule with a final decay phase near the end of training.

otherwise stated, all optimizers share the same learning rate $\eta = 3.6 \times 10^{-3}$ and weight decay $\lambda = 5 \times 10^{-3}$, with no learning-rate sweep and no optimizer-specific retuning in this appendix. All runs follow the same default seeding behavior of the training stack. For NorMuon, the most influential optimizer-specific hyperparameter is the RMS normalization, which controls scale calibration of the orthogonalized update.

A warmup–stable–decay learning-rate schedule is adopted, and results are reported for two schedule variants: (i) a warmup-enabled configuration with 10% warmup and a final 20% decay phase, and (ii) the same schedule without warmup. All other training components, including the data pipeline, batching, tokenization, model architecture, and compute budget, are kept identical across optimizers. This section is treated as a fixed-hyperparameter pilot, where broader sweeps and multi-seed evaluations are deferred to future work.

Metrics. Training loss is reported as a high-resolution indicator of early-stage optimization dynamics and stability under an identical training budget. This appendix presents the corresponding loss trajectories to complement the main-text summaries.

D.2 QWEN3-0.6B: LATE-STAGE LOSS UNDER A FIXED 1500-STEP BUDGET

The late-stage training window (steps 1200–1500) is examined to complement the early-stage analysis in the main text. Figure 2 shows the corresponding loss trajectories in this window for both schedule variants. These curves document late-stage behavior under a fixed training budget and are not intended to imply full convergence.

D.3 QWEN3-0.6B: MINIMAL REPLICATION CHECK

The main text reports Qwen3-0.6B results under a single fixed seed. To reduce the possibility that the observed qualitative trends are specific to that particular random draw, the same Qwen3-0.6B runs are repeated with an additional random seed while keeping *all* other settings identical, including the data pipeline, token budget, model architecture, learning rate, weight decay, and schedule variant. These replication runs are not used for hyperparameter tuning and are included solely as a robustness check under an identical protocol.

Figure 3 presents early-stage loss trajectories over steps 0–350 for the additional seed. Across both schedule variants, the qualitative optimizer behavior remains consistent with the main-text observations. No claim of statistical significance is made from this minimal replication; the results are provided as supplementary evidence of qualitative robustness under matched conditions.

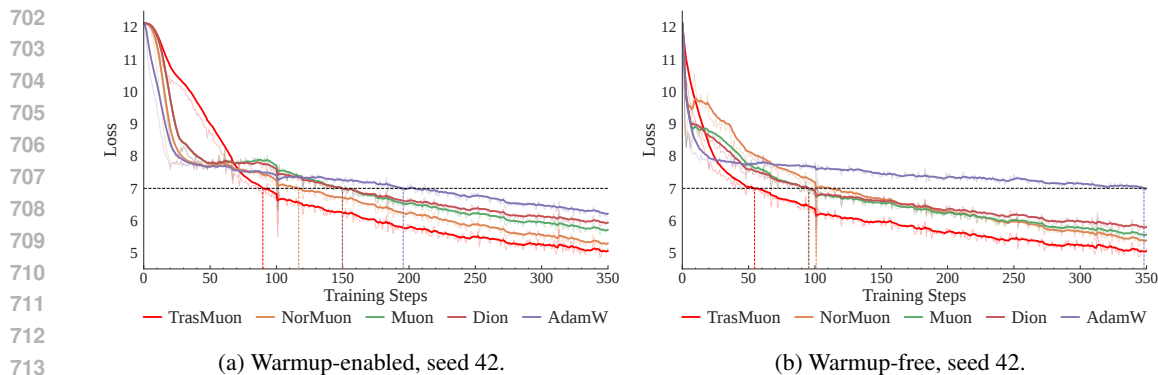


Figure 3: **Qwen3-0.6B replication under an identical protocol (additional seed)**. Early-stage training loss (steps 0–350) for an additional random seed, under the same configuration as the main-text experiment.

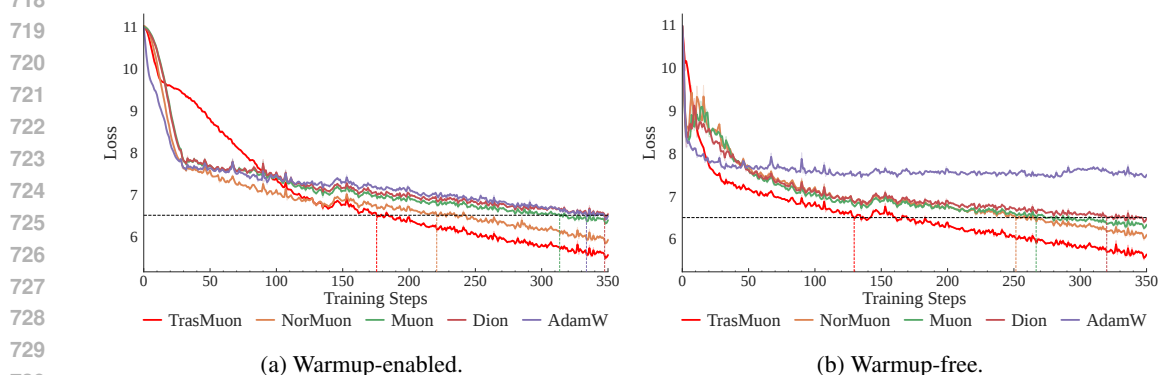


Figure 4: **GPT-2 Small early-stage dynamics (steps 0–350) under a fixed-budget protocol**. Warmup-enabled vs. warmup-free variants of the same warmup–stable–decay schedule.

D.4 GPT-2 SMALL: ADDITIONAL ARCHITECTURE UNDER THE SAME PROTOCOL

GPT-2 Small is additionally evaluated under the same fixed-budget protocol to assess whether the observed qualitative dynamics transfer to a different architecture. Figure 4 presents early-stage loss trajectories over steps 0–350 for both schedule variants, and Figure 5 reports the late-stage window over steps 1200–1500. These plots are included to complement the main-text results, and the late-stage window is reported for completeness rather than as evidence of full convergence.

E VISION TRANSFORMER EXPERIMENTS

E.1 IMAGENET-100 DATA SOURCE AND CONSTRUCTION

To reduce dataset preparation overhead, we use a publicly available ImageNet-100 *image archive* hosted on Kaggle (Ambityga, 2021) as a convenient storage source. Importantly, the Kaggle archive is used *only* as a source of image files. Class membership and the train/validation split are defined strictly according to the ILSVRC-2012 (ImageNet-1k) specification. Concretely, a fixed subset of 100 ILSVRC-2012 classes is selected (specified by synset IDs) and retained images whose labels match these synsets; we then follow the standard ILSVRC-2012 train/validation split. As a sanity check, we verify (in code) the synset-to-index mapping, per-class image counts, and that no images outside the selected synsets are included. This ensures that the benchmark corresponds to a well-defined subset of ImageNet-1k, independent of the hosting platform.

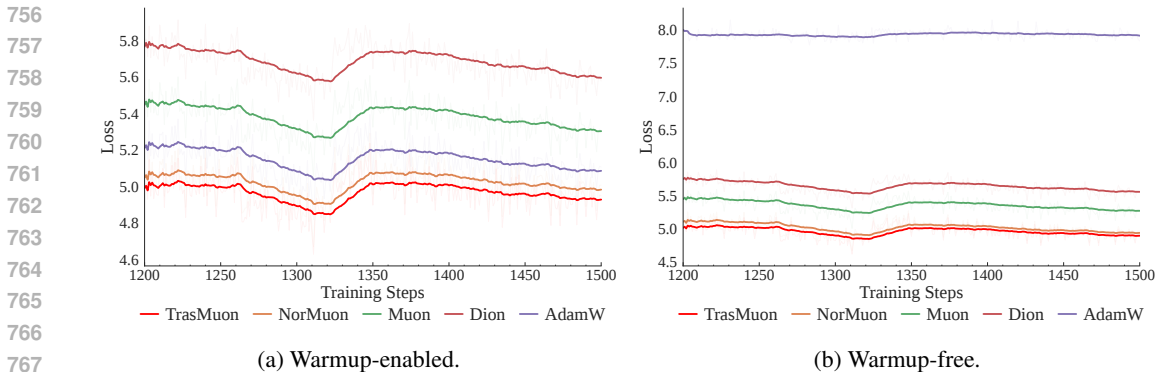


Figure 5: **GPT-2 Small late-stage window (steps 1200–1500) under a fixed-budget protocol.** Training loss for warmup-enabled vs. warmup-free schedule variants.

Table 1: **ViT on ImageNet-100.** Final *validation* top-1 accuracy (mean \pm standard deviation) over three random seeds (42, 43, 44).

Optimizer	Accuracy Mean	Accuracy Std
TrasMuon	77.47%	0.34%
NorMuon	77.10%	0.21%
Muon	69.69%	0.08%
AdamW	42.53%	4.38%

E.2 IMAGENET-100 EXPERIMENTAL PROTOCOL AND SEED ALLOCATION

We evaluate optimization methods on ImageNet-100 using a ViT-Base/16 architecture at 224×224 resolution. Training follows a standard ViT/DeiT recipe (Dosovitskiy et al., 2021; Touvron et al., 2021), including random resized cropping, horizontal flipping, color jitter, RandAugment, and random erasing, together with label smoothing and Mixup (CutMix disabled).

Weight decay is applied with parameter grouping: LayerNorm and bias parameters use zero weight decay, and all remaining parameters use a fixed decay rate. Unless otherwise stated, all experiments use a base learning rate of 1×10^{-3} and a weight decay of 5×10^{-2} . We compare AdamW (Loshchilov & Hutter, 2019), Muon (Jordan et al., 2024), NorMuon (Li et al., 2025), and TRASMUON under the same model architecture, data pipeline, training schedule, and compute budget, using three random seeds (42, 43, 44). Optimizer-specific parameters follow the respective published defaults and our unified implementation. For each method, we report the mean and standard deviation of *validation* top-1 accuracy across seeds (Table 1).

E.3 IMAGENET-100: VISION TRANSFORMER TRAINING.

We evaluate the benefits of TRASMUON in a large-scale vision setting by training ViT-Base (Dosovitskiy et al., 2021) on ImageNet-100 due to limited computational resources, a 100-class subset of ImageNet-1k (ILSVRC-2012) (Deng et al., 2009), followed the standard ILSVRC-2012 train/validation protocol. Dataset construction, class specification, and implementation details for ViT training are provided in Appendix E.1 and E.2. We compared AdamW, Muon, NorMuon, and TRASMUON, using identical training budgets and hyperparameters. Results show training loss and *validation* top-1 accuracy, aggregated over three random seeds.

Results. Across all optimizers evaluated with multi-seed runs, Muon, Normuon, and TRASMUON consistently improve optimization behavior and validation accuracy over AdamW in these experiments, demonstrating the advantage of optimizers, which update based on structured, near-orthogonal update directions. illustrated in Fig. 6a and Fig. 6b, TRASMUON achieves the fastest loss reduction, the highest validation accuracy, and reduced variability across different seeds and vari-

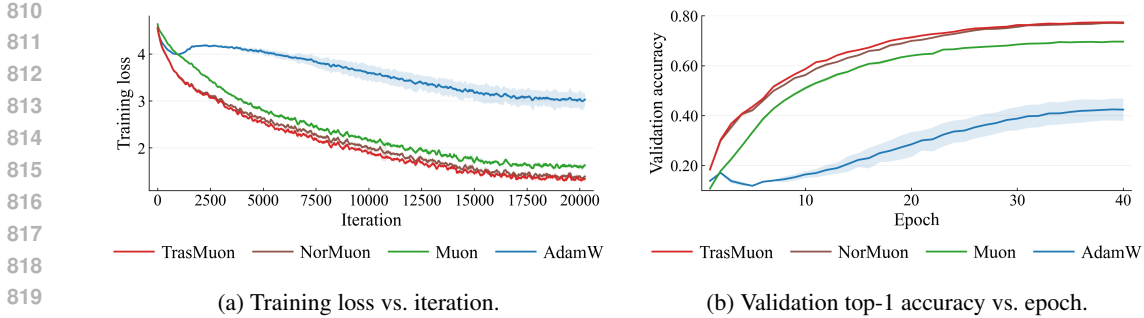


Figure 6: ViT-Base training on ImageNet-100. Multi-seed results (mean \pm std over three seeds: 42, 43, 44) for 4 optimizers. Shaded regions denote variability across seeds.

Table 2: ViT on CIFAR-100 with column-localized gradient bursts. Final top-1 test accuracy (mean \pm standard deviation) over three seeds (42, 43, 44).

Optimizer	Accuracy Mean	Accuracy Std
TrasMuon	58.77%	0.22%
NorMuon	58.31%	0.52%
Muon	57.48%	0.52%
AdamW	35.03%	5.05%

ous compared optimizers. Moreover, CIFAR-100 has been evaluated robustness under controlled column-localized burst injection in Appendix E.4.

E.4 CIFAR-100: COLUMN-ENERGY STRESS TEST

Setup (stress-test benchmark). We conduct a controlled stress test on CIFAR-100 using a Vision Transformer to assess optimizer robustness under axis-localized nonstationarity. Unless otherwise specified, all runs use 30 epochs, batch size 128, learning rate 1×10^{-3} , weight decay 5×10^{-3} , and identical data loading and preprocessing. We report mean and standard deviation of test accuracy over three random seeds (42, 43, 44).

Burst injection. To introduce structured nonstationarity without changing the data distribution, we inject sparse *column-localized gradient bursts* into selected large 2D parameter matrices (attention and MLP projections). The goal of this protocol is not to claim that injected bursts match the exact distribution of naturally occurring outliers, but to provide a reproducible mechanism-level stressor that concentrates energy along a small subset of feature axes. Crucially, the same burst pattern (targeted layers, selected columns, and timesteps) is applied across optimizers by fixing the burst random seed, enabling direct, fair comparison of optimizer responses. Full implementation details are provided in Appendix E.5.

Observed behavior. Table 2 summarizes test accuracy under burst injection. Across this stress setting, Muon-family optimizers maintain higher accuracy and lower variability than AdamW. Normalization-based variants reduce variance relative to Muon, while TRASMUON attains the highest mean accuracy and the smallest spread across seeds in this configuration.

E.5 BURST INJECTION PROTOCOL IN CIFAR-100 BENCHMARK

We define a column-wise gradient burst operator applied to selected 2D weight matrices to induce controlled column-energy spikes without altering the data distribution. At each burst step, we select k column indices (random or fixed, as specified) and perturb each selected column by adding a normalized random direction:

$$g_{:,j} \leftarrow g_{:,j} + \alpha \cdot \frac{u}{\|u\|_2 + \epsilon}, \quad u \sim \mathcal{N}(0, I).$$

864 The amplitude α can be specified as a fixed absolute value or scaled relative to the current gradient
865 magnitude using a Frobenius-normalized reference:

$$866 \alpha = \rho \cdot \frac{\|g\|_F}{\sqrt{d_{\text{out}}d_{\text{in}}}},$$

869 optionally clipped by a maximum threshold. Bursts occur every T optimization steps after an op-
870 tional warmup phase and target only designated 2D layers (e.g., attention projections and MLP
871 weights). Burst events and optimizer internal statistics (including feature-wise clipping coefficients,
872 when applicable) are logged at the same timesteps to enable direct alignment between perturbations
873 and optimizer responses.

874 F PINNS WITH RANDOM-ROI SAMPLING STRESS TEST

875 Adaptive collocation in physics-informed neural networks (PINNs) is often necessary to resolve lo-
876 calized errors and stiff PDE behavior, where uniform sampling under-resolves difficult regions (Gao
877 et al., 2023; Subramanian et al., 2022; Wu et al., 2023). Here we use region-of-interest (ROI) densi-
878 fication as a *controlled nonstationarity* mechanism to stress-test optimizer robustness: periodically
879 concentrating interior collocation points in a small subregion induces distribution shifts in the resid-
880 ual samples, perturbing gradient statistics in a reproducible way.

883 **Helmholtz equation setup.** On $\Omega = [0, 1]^2$, we consider

$$884 \Delta u(x) + \kappa^2 u(x) = f(x), \quad u|_{\partial\Omega} = 0, \quad (38)$$

885 with the manufactured solution

$$886 u^*(x, y) = \sin(\pi k x) \sin(\pi k y), \quad \kappa = \pi k. \quad (39)$$

887 which yields $f(x) = -(\pi k)^2 u^*(x)$.

888 We train an MLP $u_\theta : \Omega \rightarrow \mathbb{R}$ by minimizing

$$889 \mathcal{L}(\theta) = \mathbb{E} \left[\frac{1}{2} r_\theta(x)^2 \right] + \lambda_b \mathbb{E} \left[\frac{1}{2} (u_\theta(x) - u^*(x))^2 \right] \quad (40)$$

$$891 r_\theta(x) = \Delta u_\theta(x) + \kappa^2 u_\theta(x) - f(x). \quad (41)$$

892 and report the relative error on a fixed evaluation grid

$$893 \text{rel-}L_2(u_\theta, u^*) = \frac{\|u_\theta - u^*\|_2}{\|u^*\|_2}. \quad (42)$$

894 **ROI sampling protocol.** To emulate adaptive densification, we impose a controlled, time-varying
895 interior sampling distribution. We consider the Helmholtz equation on $\Omega = [0, 1]^2$ with homo-
896 geneous Dirichlet boundary conditions and a manufactured solution u^* . We run 4000 optimiza-
897 tion steps and introduce nonstationary ROI events after step $t_0 = 1000$. At each step we sample
898 $N_r = 1024$ interior points and $N_b = 256$ boundary points, with boundary weight $\lambda_b = 100$. We
899 evaluate every 200 steps on a fixed 128×128 grid.

900 **Nonstationary ROI events (distribution shift).** Starting from step $t_0 = 1000$, we trigger ROI
901 events every $K_{\text{out}} = 20$ steps. At an ROI event step t , interior points are sampled from a mixture
902 distribution

$$903 p_t(x) = (1 - \alpha) p_0(x) + \alpha p_{\text{roi}}^{(t)}(x), \quad \alpha = 0.95, \quad (43)$$

904 where p_0 is uniform over Ω and $p_{\text{roi}}^{(t)}$ is uniform over the selected ROI patch $\Omega_{\text{roi}}^{(t)} = [x_0, x_1] \times [y_0, y_1]$.
905 This yields repeated, time-varying distribution shifts that mimic practical ROI/adaptive refinement
906 policies in PINNs.

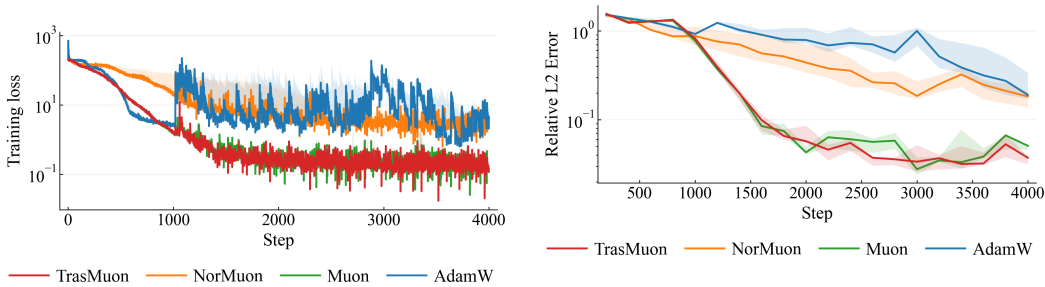
907 **Random ROI patch pool (reproducible).** To avoid sensitivity to a single ROI location, $\Omega_{\text{roi}}^{(t)}$ is
908 chosen from a fixed pool (Table 3) using a deterministic `step_hash` seeding rule. Thus ROI
909 locations vary across events while remaining fully reproducible given the experiment configuration
910 and the training step index.

Table 3: ROI patch pool used for ROI events (rectangles are $[x_0, x_1] \times [y_0, y_1]$).

ROI patches
Corners: $[0.00, 0.03] \times [0.00, 0.03]$, $[0.97, 1.00] \times [0.00, 0.03]$, $[0.00, 0.03] \times [0.97, 1.00]$, $[0.97, 1.00] \times [0.97, 1.00]$
Edges: $[0.48, 0.53] \times [0.00, 0.05]$, $[0.48, 0.53] \times [0.95, 1.00]$, $[0.00, 0.05] \times [0.48, 0.53]$, $[0.95, 1.00] \times [0.48, 0.53]$
Interior: $[0.20, 0.25] \times [0.20, 0.25]$, $[0.45, 0.50] \times [0.10, 0.15]$, $[0.10, 0.15] \times [0.55, 0.60]$, $[0.60, 0.65] \times [0.60, 0.65]$

ROI-local evaluation. To quantify localized disturbance and recovery, at ROI event steps we additionally compute an ROI-local $\text{rel}L_2$ on a 64×64 grid restricted to $\Omega_{\text{roi}}^{(t)}$. We estimate non-ROI error by sampling 4096 points from $\Omega \setminus \Omega_{\text{roi}}^{(t)}$. These diagnostics separate global convergence from localized behavior under distribution shifts.

F.1 PINNS BENCHMARK: ROI SAMPLING AS A NONSTATIONARY STRESS TEST



(a) training objective (estimated on the time-varying sampling distribution p_t)

(b) domain-wide relative L_2 error evaluated on a fixed grid. ROI events start at step 1000 and repeat every 20 steps

Figure 7: PINN Helmholtz ($k=2$) under random ROI sampling shifts. Curves show the mean over seeds, and shaded regions indicate variability across seeds.

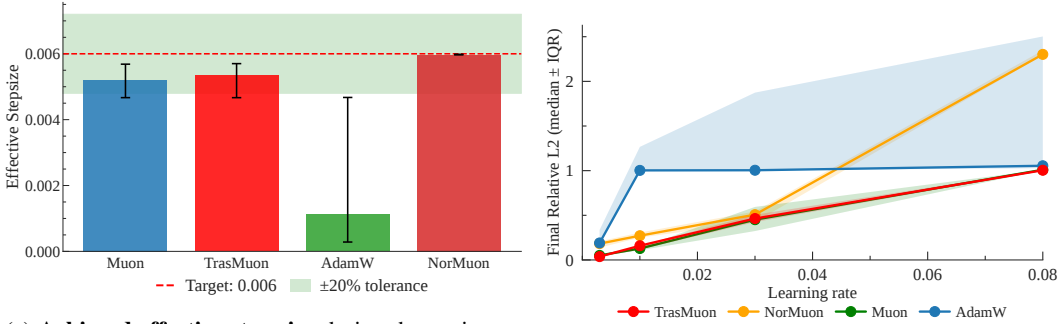
PINN ROI-sampling stress test: convergence and robustness. Figure 7a and 7b compare Muon and TRASMUON on Helmholtz ($k=2$) under a controlled nonstationary ROI-sampling protocol, where ROI events start at step 1000 and recur every 20 steps. During the initial stationary phase (before ROI events), both methods exhibit nearly identical optimization trajectories in terms of training loss and domain-wide relative L_2 error, indicating that TRASMUON does not incur a measurable overhead or degradation under standard uniform sampling.

After ROI events begin, the training objective becomes significantly more variable due to the induced distribution shifts in interior collocation points. In this nonstationary regime, TRASMUON maintains comparable or slightly lower training loss while exhibiting reduced extreme fluctuations, consistent with its design goal of suppressing bursty, feature-localized updates. These results support the conclusion that TRASMUON preserves baseline convergence under stationary sampling, while improving stability and final solution accuracy under controlled, nonstationary ROI sampling shifts.

F.2 STEP-SIZE ALIGNMENT AND LEARNING-RATE SENSITIVITY (PINNS HELMHOLTZ, $k = 2$)

Motivation. Nominal learning rates are not directly comparable across update rules because different optimizers can induce different *effective* parameter-space step magnitudes. To reduce the confound that performance differences are driven by trivial step-size mismatches, we complement the main comparison with (i) a short step-size alignment diagnostic and (ii) a shared learning-rate (LR) sweep.

Step-size alignment diagnostic. Figure 8a reports the achieved effective step size, computed from parameter differences during an initial stationary window (before any ROI perturbations are intro-



(a) **Achieved effective step size** during the stationary alignment window (before ROI events). The dashed line is the target, and the shaded region indicates tolerance. Error bars denote variability across seeds. (b) **LR sensitivity** of final relative L_2 error under a shared LR sweep. Lines show the median across seeds and shaded bands indicate the IQR.

Figure 8: **PINNs Helmholtz ($k = 2$): step-size alignment and LR sensitivity.** (a) Effective step-size alignment reduces trivial magnitude confounds when comparing optimizers with different update rules. (b) A shared LR sweep summarizes sensitivity to step-size calibration via final relative L_2 error.

duced),

$$s_t = \frac{\|\Delta\theta_t\|_2}{\sqrt{P}}, \quad P = \dim(\theta).$$

We target a fixed reference magnitude (dashed line) with a tolerance band (shaded region). Muon and TRASMUON attain comparable achieved step sizes within the tolerance range, supporting that subsequent robustness comparisons are not explained by a simple global step-size discrepancy.

Learning-rate sweep. Figure 8b shows the final relative L_2 error under a shared LR sweep. Shaded bands summarize variability across random seeds (median with an interquartile range). Both methods exhibit the expected degradation as LR increases beyond the stable region. Together with the alignment diagnostic, this sweep provides a complementary view of optimizer sensitivity to step-size calibration under the same training and sampling protocol.

F.3 PINNS DIAGNOSTICS: METRIC DISTRIBUTIONS ACROSS SEEDS

We visualize the distribution of key robustness and accuracy metrics across random seeds for the PINN Helmholtz benchmark ($k = 2$). This figure serves as a distributional check to ensure that the reported trends are not driven by a single favorable run.

G CONTROLLED DIAGNOSTICS

This appendix provides protocol-level details and supporting evidence for the controlled diagnostics study. The intent is two-fold: (i) to make the stress protocol fully reproducible, and (ii) to document a minimal, time-aligned evidence chain that is *consistent with* the intended energy-indexed, feature-wise clipping mechanism under a controlled intervention. We do not introduce new claims beyond Section G.1.

G.1 CONTROLLED DIAGNOSTICS: COLUMN-LOCALIZED OUTLIERS AND ENERGY-BASED FEATURE CLIPPING

We design a controlled toy problem to validate the *feature-wise clipping* mechanism in TRASMUON under intermittent, column-localized bursts. We optimize a matrix parameter $W \in \mathbb{R}^{d \times d}$ under the quadratic objective

$$\min_W f(W) = \frac{1}{2} \|AWB - T\|_F^2, \tag{44}$$

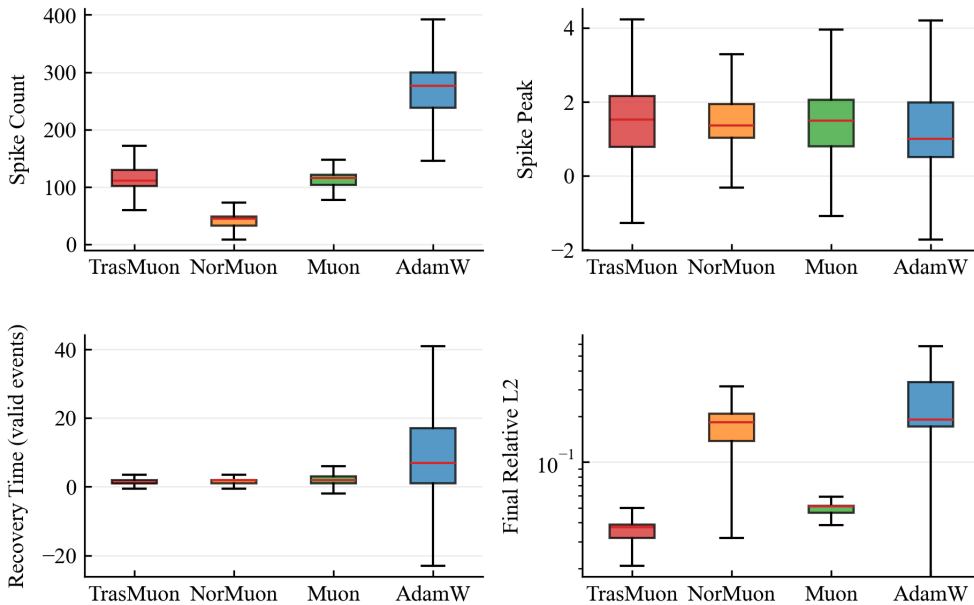


Figure 9: **PINNs Helmholtz($k = 2$): metric distributions across seeds.** Boxplots summarize spike count, spike peak, recovery time (valid events), valid-event rate, and final relative L_2 error across seeds under the same training protocol. Each box shows the median and interquartile range (IQR); whiskers indicate the remaining spread.

where $A = U\Sigma_A U^\top$ and $B = V\Sigma_B V^\top$ with random orthogonal U, V , and diagonal spectra Σ_A, Σ_B chosen to yield a target condition number $\kappa \in \{10^2, 10^4, 10^6\}$. This construction controls stiffness while allowing us to randomize nuisance rotations across runs.

Column-localized outlier injection (stress protocol). To emulate rare, feature-localized gradient domination, every K_{out} steps we inject an outlier event that amplifies a small subset of columns in a fixed feature basis. Concretely, for momentum M_t we select a set \mathcal{J} of $s \ll d$ column indices and apply a multiplicative burst

$$\widetilde{M}_{t,j} = \begin{cases} a M_{t,j}, & \text{if } j \in \mathcal{J}, \\ M_{t,j}, & \text{otherwise,} \end{cases} \quad (45)$$

with burst amplitude $a > 1$. This perturbation produces abrupt increases in column energy $E_{t,j} = \sum_i \widetilde{M}_{t,ij}^2$ while leaving the underlying objective equation 44 unchanged.

Preserving feature semantics. Because TRASMUON’s clipping is axis-aligned (column-wise), the stress protocol is evaluated under a `fix_V=True` setting, i.e., the column basis is preserved across training and across injected events. We additionally report a boundary condition where the column basis is randomized (`fix_V=False`); in that case, injected energy disperses across columns and feature-wise clipping is not expected to yield an advantage.

Metrics. We track (i) spike count and (ii) final objective value, reporting median and IQR over multiple seeds/rotations.

Closed-loop response. Figure 10 shows that TRASMUON reduces burst-induced loss spikes and improves convergence relative to the NorMuon backbone under matched compute. Figure 11 provides mechanism-level evidence consistent with a closed-loop response: outlier events increase the relative column-energy ratio (e.g., r_{q95}/r_{max}), which is immediately followed by a decrease in the *applied* clipping signal (tracked by $c_{\text{used},\text{min}}$), thereby damping the burst and suppressing spikes.

1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

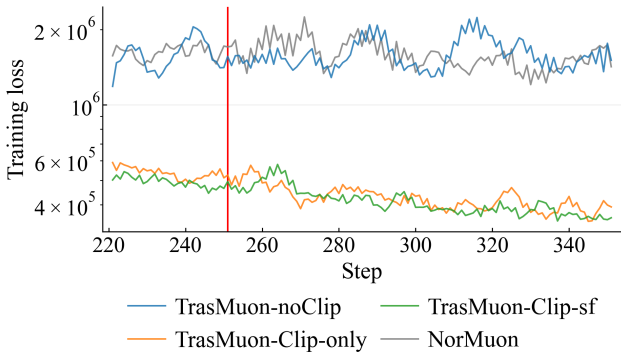


Figure 10: **Column outlier injection.** Loss trajectories in a window around an outlier event. Vertical markers indicate outlier steps.

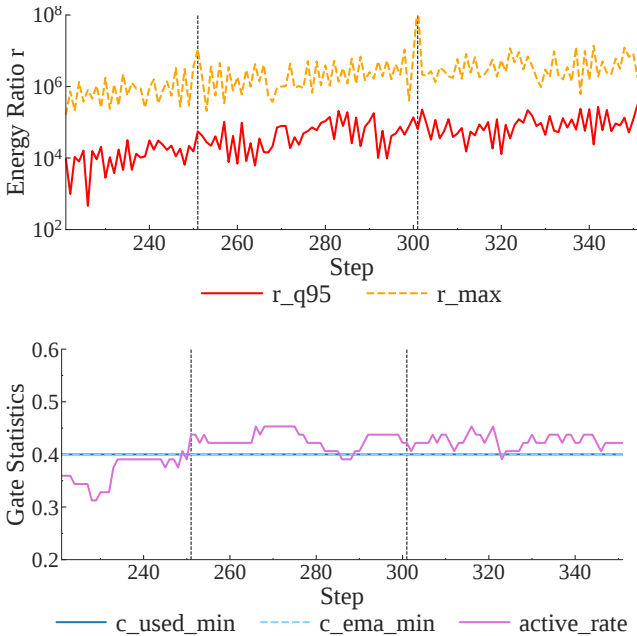


Figure 11: **Closed-loop clipping evidence.** Outlier events increase the column-energy ratio in log-scale (top), followed by stronger *feature-wise clipping* in the applied coefficients (bottom; $c_{used,min}$).

Not a trivial step-size reduction. To rule out the confound that improvements arise from a global effective step-size change, we include a TRASMUON-NOCLIP ablation that disables feature-wise clipping while keeping all other components identical. As summarized in Table 4, TRASMUON-NOCLIP behaves similarly to NorMuon, whereas enabling clipping yields a clear reduction in spike statistics and a large improvement in the final objective, isolating the contribution of feature-wise clipping.

Boundary condition (feature semantics broken). When the column basis is randomized, the advantage of feature-wise clipping diminishes, consistent with the intended mechanism: axis-aligned clipping requires a meaningful feature basis.

G.2 MINIMAL EVIDENCE CHAIN (TIME-ALIGNED OBSERVABLES)

Section G.1 argues that spike suppression is consistent with the following ordering under a controlled intervention: *outlier injection* \rightarrow *relative energy increases* \rightarrow *stronger applied clipping* \rightarrow *reduced*

Table 4: **Controlled diagnose summary under `fix_V=True`.** TRASMUON-NOCLIP removes feature-wise clipping while keeping the rest identical. We report median with IQR across runs; lower is better.

Method	Spike Count	Final Loss
NorMuon	44 (35,56)	1.3e+06 (1.0e+06,1.6e+06)
TRASMUON-noClip	48 (38,56)	1.1e+06 (8.9e+05,1.9e+06)
TRASMUON-Clip-only	28 (24,34)	2.4e+05 (2.0e+05,2.8e+05)
TRASMUON-Clip-sf	30 (24,36)	2.0e+05 (1.6e+05,2.7e+05)

loss spikes. We summarize the corresponding observables, which are directly logged and visualized in Fig. 11.

(1) Outlier injection increases relative energy. Eq. equation 45 increases the column energies $E_{t,j} = \sum_i \widetilde{M}_{t,i,j}^2$ for $j \in \mathcal{J}$. We track the relative energy ratio

$$r_{t,j} = \frac{E_{t,j}}{E_t^{\text{ref}} + \epsilon}, \quad (46)$$

where E_t^{ref} is the running reference used by TRASMUON. We visualize robust summaries such as r_{q95} and r_{\max} , which rise at injected outlier steps (Fig. 11, top).

(2) Higher relative energy is followed by stronger applied clipping. TRASMUON produces damping-only clipping coefficients $c_{t,j}^{\text{used}} \in [c_{\min}, 1]$ that decrease with $r_{t,j}$. Consistent with this design, outlier steps are followed by a decrease in the applied signal, visible via $c_{\text{used},\min} = \min_j c_{t,j}^{\text{used}}$ (Fig. 11, bottom). This time alignment is consistent with the intended ordering “energy rise \rightarrow clipping increase” under the intervention.

(3) Applied clipping attenuates column updates (selectively). Given the matrix-form update,

$$\Delta W_t = -\hat{\eta}_t O_t^{\text{base}} \text{diag}(c_t^{\text{used}}), \quad (47)$$

each column update magnitude is scaled by $c_{t,j}^{\text{used}}$:

$$\|\Delta W_{t,j}\|_2 = c_{t,j}^{\text{used}} \hat{\eta}_t \|O_{t,j}^{\text{base}}\|_2 \leq \hat{\eta}_t \|O_{t,j}^{\text{base}}\|_2. \quad (48)$$

Thus clipping is *selective*: only columns with $c_{t,j}^{\text{used}} < 1$ are damped, while the structured direction O_t^{base} is preserved.

(4) Spike suppression and objective improvement. Consistent with (1)–(3), TRASMUON reduces loss spikes around outlier events (Fig. 10) and achieves lower final objective values than the backbone under matched compute (Table 4). Spike metrics (count/peak) are computed using the same deterministic detection rule across methods; details are provided in our experiment scripts and plotting utilities.

Table 5: **Boundary condition effects(`fix_V=False`).** When feature/column semantics are broken by column-space mixing, the advantage of feature-wise clipping diminishes.

Method	Spike Count	Final Loss
NorMuon	79 (74,86)	1.1e+06 (8.4e+05,1.6e+06)
TRASMUON-noClip	80 (74,87)	1.4e+06 (9.8e+05,1.7e+06)
TRASMUON-clip-only	74 (67,80)	1.5e+06 (1.2e+06,1.8e+06)
TRASMUON-clip+SF	72 (65,80)	1.3e+06 (9.4e+05,1.7e+06)

Controls that break the chain. We include two controls that remove key requirements of the mechanism: (i) TRASMUON-NOCLIP sets $c_t^{\text{used}} \equiv 1$, removing the attenuation in Eq. equation 48; empirically it behaves similarly to NorMuon (Table 4). (ii) Under `fix_V=False`, the injected energy is dispersed across columns, so clipping is no longer aligned with injected directions; correspondingly, the advantage of feature-wise clipping diminishes (Table 5).