

How Does Orthogonalization Adapt to the Neural-Network Hessian Structure? A Gradient Self Outer-Product Analysis at Initialization

author names withheld

Under Review for the Workshop on High-dimensional Learning Dynamics, 2026

Abstract

Muon orthogonalizes a weight matrix’s momentum before each step, and on neural networks this simple preconditioner beats entry-wise optimizers by a wide margin. Most existing analyses, however, work in a very abstract problem class, from which it is hard to see why orthogonalization should be particularly suited to neural networks. This work analyzes Muon’s preconditioner in three concrete neural-network settings. The layer-wise Hessian of a neural network is known to be diagonally dominant within its row blocks. While Muon’s implicit preconditioner has the matching Kronecker form $(\mathbf{V}\mathbf{V}^\top)^{1/2} \otimes \mathbf{I}$. The two align exactly when $\mathbf{V}\mathbf{V}^\top$ is itself diagonal, which raises a concrete question: when is $\mathbf{V}\mathbf{V}^\top$ (approximately) diagonal? To answer it, we compute $\mathbb{E}[\mathbf{G}\mathbf{G}^\top]$ (equivalently $\mathbb{E}[\mathbf{V}\mathbf{V}^\top]$ at initialization) in closed form under Gaussian init, for three standard settings: symmetric matrix factorization, deep linear networks, and two-layer ReLU networks. In each case, the diagonal entries dominate the off-diagonal ones as the width grows. Hence $\mathbf{V}\mathbf{V}^\top$ is asymptotically diagonal: Muon’s preconditioner aligns with the Hessian’s row-block structure.

1. Introduction

Muon [14] optimizes a matrix-shaped weight $\mathbf{W}_t \in \mathbb{R}^{m \times n}$ by orthogonalizing the momentum buffer instead of using it directly. From a stochastic gradient \mathbf{G}_t ,

$$\mathbf{V}_t = \beta \mathbf{V}_{t-1} + (1 - \beta) \mathbf{G}_t, \quad \mathbf{D}_t = \text{NS}_5(\mathbf{V}_t) \approx (\mathbf{V}_t \mathbf{V}_t^\top)^{-1/2} \mathbf{V}_t, \quad \mathbf{W}_{t+1} = \mathbf{W}_t - \eta_t \mathbf{D}_t, \quad (1)$$

with $\mathbf{V}_0 = \mathbf{0}$ and NS_5 a fixed-point iteration for the orthogonal polar factor: if $\mathbf{V}_t = \mathbf{P}\Sigma\mathbf{Q}^\top$ is a thin SVD then $\mathbf{D}_t = \mathbf{P}\mathbf{Q}^\top = (\mathbf{V}_t \mathbf{V}_t^\top)^{-1/2} \mathbf{V}_t$. A growing body of work studies Muon as a general optimization algorithm: momentum steepest descent, or a non-Euclidean trust-region step, in a spectral-norm geometry, with convergence rates and an implicit spectral-norm constraint on the iterates [1–4, 6, 15–17, 21, 24].

All of this works in an abstract problem class, and that is exactly its limitation: it cannot explain why Muon, a simple matrix preconditioner, dramatically outperforms entry-wise optimizers *on neural networks specifically*. The few network-specific results [9, 19, 26–28, 32] sharpen the picture in particular regimes but leave the basic structural question open: what makes the orthogonalization preconditioner suited to a neural network’s curvature?

Preconditioner shape versus Hessian shape. Take a single step from initialization, where $\mathbf{V}_0 = \mathbf{G}_0$ and hence $\mathbf{V}_t \mathbf{V}_t^\top = \mathbf{G}_t \mathbf{G}_t^\top$. Vectorizing (1) row by row gives $\text{vec}(\mathbf{D}_t) = ((\mathbf{V}_t \mathbf{V}_t^\top)^{-1/2} \otimes \mathbf{I}_n) \text{vec}(\mathbf{V}_t)$, so Muon descends in the metric

$$H_{\text{Muon}} = \mathbf{P}_t^{1/2} \otimes \mathbf{I}_n, \quad \mathbf{P}_t := \mathbf{V}_t \mathbf{V}_t^\top \in \mathbb{R}^{m \times m}, \quad (2)$$

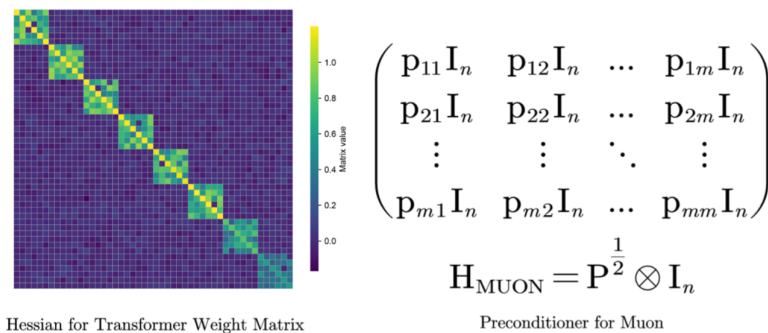


Figure 1: Muon’s implicit preconditioner $H_{\text{MUON}} = \mathbf{P}_t^{1/2} \otimes \mathbf{I}_n$ next to the layer-wise Hessian. When $\mathbf{P}_t = \mathbf{V}_t \mathbf{V}_t^\top$ is close to diagonal, H_{MUON} matches the dominant row-block pattern of the neural-network Hessian.

a Kronecker shape determined entirely by \mathbf{P}_t . On the other side, the layer-wise Hessian of a weight $\mathbf{W} \in \mathbb{R}^{m \times n}$, unfolded into $m \times m$ blocks of size $n \times n$ along the m row directions, is *row-block diagonally dominant* [8, 33, 34]; since the Hessian is the ideal preconditioner under a quadratic model, that block pattern is the shape a good preconditioner should have. By (2), H_{MUON} takes exactly that block-diagonal form precisely when \mathbf{P}_t is diagonal (Figure 1). So the question reduces to a clean property of the gradient.

Question. Is $\mathbf{V}_t \mathbf{V}_t^\top$ (equivalently, $\mathbf{G}_t \mathbf{G}_t^\top$ at initialization) diagonally dominant for neural-network weight matrices?

For neural networks the answer is yes, in the limit of large width. In practice this has been verified empirically along the entire training trajectory of GPT-2 and LLaMA [7]; here we give a theoretical answer at initialization in the three settings most used in neural-network theory, with the diagonal-versus-off-diagonal orders summarized at the end of Section 4. Two facts come out of the analysis.

1. The orthogonalization preconditioner aligns asymptotically with the neural-network Hessian: $\mathbf{P}_t = \mathbf{V}_t \mathbf{V}_t^\top$ becomes diagonal as the width grows, so H_{MUON} becomes block-diagonal in the Hessian’s row-block structure.
2. Orthogonalization is asymptotically the row-wise normalization used by several recent stateless optimizers: when $\mathbf{V}_t \mathbf{V}_t^\top$ is diagonal, $(\mathbf{V}_t \mathbf{V}_t^\top)^{-1/2} \mathbf{V}_t$ is exactly each row of \mathbf{V}_t divided by its own ℓ_2 norm. (Appendix D)

2. Related Work

Theory of Muon. One line analysis Muon in a general optimization setup (such as nonconvex smooth). Li and Hong [17], Shen et al. [24] give non-convex convergence rates for momentum

steepest descent in the spectral norm and its practical variants; Kovalev [15] reads gradient orthogonalization as a non-Euclidean trust-region step; Bernstein and Newhouse [2, 3], Chen et al. [4], Pethick et al. [21] describe Muon through norm-constrained linear minimization oracles and modular duality, with Chen et al. [4] extracting an implicit spectral-norm constraint on the weights; An et al. [1], Lau et al. [16] fit Muon into a wider family of matrix-gradient preconditioners; and Davis and Drusvyatskiy [6] asks when spectral updates help. None of this depends on the network, so none of it addresses why orthogonalization matches neural-network curvature. A smaller line does look at the network: Fan et al. [9], Vasudeva et al. [27] characterize the margin bias of spectral descent and its generalization benefit on separable or imbalanced data; Wang et al. [28] studies Muon on associative-memory learning; Zhang et al. [32] ties layer-wise preconditioning to provable feature learning; Su [26] asks whether orthogonalization is optimal under an isotropic curvature model; and Ma et al. [19] raises, in a matrix-factorization setup, the very question of how gradient orthogonalization meets the curvature of training. We answer it through the diagonal structure of $\mathbb{E}[\mathbf{G}\mathbf{G}^\top]$ across the three settings.

Row-wise normalization optimizers. A separate line drops optimizer state by normalizing the gradient row by row. Zhang et al. [34] groups learning rates by Hessian block, motivated directly by the row-block structure of the Transformer Hessian; Glentis et al. [10], Ma et al. [18], Scetbon et al. [23], Wen et al. [29] replace Adam’s state with a row-wise ℓ_2 normalization (plus a second, cheap normalization); Gu and Xie [11], Xu et al. [31] look at the manifold-optimization and width-scaling sides of row/column normalization; Pethick et al. [21] reach a row-wise normalization from norm-constrained steepest descent; and Deng et al. [7] make the link to Muon explicit, showing row-momentum normalization matches Muon at scale. Our analysis explains the common thread: once $\mathbf{V}_t\mathbf{V}_t^\top$ is diagonally dominant, row normalization and orthogonalization are the same operation.

3. Preliminaries

We study, at a Gaussian initialization, the $m \times m$ matrix $\mathbb{E}[\mathbf{G}\mathbf{G}^\top]$ for a weight $\mathbf{W} \in \mathbb{R}^{m \times n}$ with population loss \mathcal{L} and gradient $\mathbf{G} = \partial\mathcal{L}/\partial\mathbf{W}$; this is the expected value of $\mathbf{P}_0 = \mathbf{V}_0\mathbf{V}_0^\top$ in (2). We use three settings, each stated with its own network and assumption. Throughout, \mathbf{I}_d is the identity, tr the trace, $\|\cdot\|_F$ the Frobenius norm, $\sigma(z) = \max(z, 0)$ the ReLU with $\sigma'(z) = \mathbf{1}[z > 0]$, \otimes the Kronecker product, and $a_n = \Theta(b_n)$ means $c b_n \leq a_n \leq C b_n$ for fixed $0 < c \leq C$ as the relevant dimension grows.

Symmetric matrix factorization. For positive integers d, k, r with $k \geq r$, the over-parameterized loss is $\mathcal{L}(\mathbf{U}) := \frac{1}{4} \|\mathbf{U}\mathbf{U}^\top - \mathbf{M}^*\|_F^2$ with $\mathbf{U} \in \mathbb{R}^{d \times k}$ and a fixed target $\mathbf{M}^* \in \mathbb{R}^{d \times d}$ that is symmetric positive semidefinite of rank r ; the gradient is $\mathbf{G} = (\mathbf{U}\mathbf{U}^\top - \mathbf{M}^*)\mathbf{U} \in \mathbb{R}^{d \times k}$. The width is k .

Deep linear network. For $L \geq 1$ and layer dimensions d_0, \dots, d_L , the network $f(\mathbf{x}) = \mathbf{W}_L \cdots \mathbf{W}_1 \mathbf{x}$ with $\mathbf{W}_i \in \mathbb{R}^{d_i \times d_{i-1}}$ learns a fixed linear target $\Phi \in \mathbb{R}^{d_L \times d_0}$ under $\mathcal{L} = \frac{1}{2} \mathbb{E}_{\mathbf{x}} \|f(\mathbf{x}) - \Phi \mathbf{x}\|^2$. Write $\mathbf{B}_i := \mathbf{W}_{i-1} \cdots \mathbf{W}_1$ (with $\mathbf{B}_1 := \mathbf{I}$), $\mathbf{A}_i := \mathbf{W}_L \cdots \mathbf{W}_{i+1}$ (with $\mathbf{A}_L := \mathbf{I}$), and $\mathbf{R} := \mathbf{W}_L \cdots \mathbf{W}_1 - \Phi$; the layer- i gradient is $\mathbf{G}_i = \partial\mathcal{L}/\partial\mathbf{W}_i = \mathbf{A}_i^\top \mathbf{R} \mathbf{B}_i^\top$. The widths are the inner dimensions d_1, \dots, d_{L-1} .

Two-layer ReLU network. For positive d_0, d_1, d_2 , the network $f(x) = W_2\sigma(W_1x)$ with $W_1 \in \mathbb{R}^{d_1 \times d_0}$, $W_2 \in \mathbb{R}^{d_2 \times d_1}$ learns a fixed linear target $\Phi \in \mathbb{R}^{d_2 \times d_0}$ under $\mathcal{L} = \frac{1}{2}\mathbb{E}_x \|f(x) - \Phi x\|^2$; write $h(x) := \sigma(W_1x)$, w_a for the a -th row of W_1 , and $G_i := \partial\mathcal{L}/\partial W_i$. The hidden dimension is d_1 .

Assumption 1 (Gaussian initialization) *All weight entries (\mathbf{U} ; every \mathbf{W}_i ; W_1 and W_2) are i.i.d. $\mathcal{N}(0, 1)$; in the deep linear and ReLU settings the input is $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_0})$ and is integrated out in \mathcal{L} ; the target (\mathbf{M}^* or Φ) is deterministic and independent of the weights. Expectations are over the weights.*

4. Main Results

We state the closed form of $\mathbb{E}[\mathbf{G}\mathbf{G}^\top]$ in each setup, read off the dominance, and defer the proofs to Sections A to C.

Symmetric matrix factorization.

Theorem 1 *Under Assumption 1 (matrix factorization), with $\mathbf{W} := \mathbf{U}\mathbf{U}^\top \sim \mathcal{W}_d(k, \mathbf{I}_d)$ Wishart,*

$$\mathbb{E}[\mathbf{G}\mathbf{G}^\top] = \alpha_3 \mathbf{I}_d - 2\alpha_2 \mathbf{M}^* + k (\mathbf{M}^*)^2, \quad (3)$$

where $\alpha_2 = k(k + d + 1)$ and $\alpha_3 = k^3 + 3(d + 1)k^2 + (d^2 + 3d + 4)k$ are the (isotropic) diagonal entries of $\mathbb{E}[\mathbf{W}^2]$ and $\mathbb{E}[\mathbf{W}^3]$.

Reading off entries with d, \mathbf{M}^* fixed: the diagonal is $\alpha_3 - 2\alpha_2 \mathbf{M}_{ii}^* + k(\mathbf{M}_{ii}^*)^2 = \Theta(k^3)$ and the off-diagonal is $-2\alpha_2 \mathbf{M}_{ij}^* + k(\mathbf{M}_{ij}^*)^2 = \Theta(k^2)$ whenever $\mathbf{M}_{ij}^* \neq 0$, so the ratio is $\Theta(k/\mathbf{M}_{ij}^*)$. As the width k grows, $\mathbb{E}[\mathbf{G}\mathbf{G}^\top]$ is diagonal up to a factor k , and the large isotropic $\Theta(k^3)$ part is exactly what orthogonalization strips away, leaving the target-aligned \mathbf{M}^* .

Deep linear network.

Theorem 2 *Under Assumption 1 (deep linear network), for every layer $i \in \{1, \dots, L\}$,*

$$\mathbb{E}[\mathbf{G}_i \mathbf{G}_i^\top] = V_i s_i \mathbf{I}_{d_i} + \mathbf{T}_i, \quad \mathbf{T}_i = \begin{cases} (\prod_{j=1}^{L-1} d_j) \Phi \Phi^\top, & i = L, \\ (\prod_{j \neq i, 1 \leq j \leq L-1} d_j) \|\Phi\|_{\mathbb{F}}^2 \mathbf{I}_{d_i}, & 1 \leq i < L, \end{cases} \quad (4)$$

where V_i, s_i are positive dimension-only scalars given by two-step recursions in (d_0, \dots, d_L) (Section B), with $V_L s_L \asymp d_{L-1} (\prod_{j=0}^{L-2} d_j)^2$.

For a hidden layer ($i < L$) the off-diagonal block \mathbf{T}_i is a multiple of \mathbf{I}_{d_i} , so $\mathbb{E}[\mathbf{G}_i \mathbf{G}_i^\top]$ is exactly diagonal: linearity together with the sign symmetry $\mathbf{W}_i \mapsto \mathbf{D}\mathbf{W}_i$, $\mathbf{W}_{i+1} \mapsto \mathbf{W}_{i+1}\mathbf{D}^{-1}$ for diagonal \mathbf{D} kills every off-diagonal expectation. At the output layer the diagonal noise $V_L s_L$ grows with the inner widths while the off-diagonal stays proportional to $(\prod_{j < L} d_j) (\Phi \Phi^\top)_{ab}$, so the ratio still diverges.

Two-layer ReLU network. Let $\theta := \angle(u, v)$ for independent $u, v \sim \mathcal{N}(0, \mathbf{I}_{d_0})$, and set the dimension-free constants $\kappa_1 := \mathbb{E}\left[\frac{\sin^2 \theta + (\pi - \theta)^2}{4\pi^2} + \frac{\sin \theta \cos \theta (\pi - \theta)}{2\pi^2}\right]$ and $\kappa_2 := \frac{\mathbb{E}[\sin \theta (\pi - \theta)]}{2\pi^2}$; as $d_0 \rightarrow \infty$, $\kappa_1 \rightarrow \frac{1}{4\pi^2} + \frac{1}{16}$ and $\kappa_2 \rightarrow \frac{1}{4\pi}$.

Theorem 3 *Under Assumption 1 (two-layer ReLU network), with $H := \mathbb{E}_x[h(x)h(x)^\top]$ and $\alpha := \mathbb{E}_{W_1}[\text{tr}(H^2)]$,*

$$\mathbb{E}[G_2 G_2^\top] = \alpha \mathbf{I}_{d_2} + \frac{d_1}{4} \Phi \Phi^\top, \quad \alpha = \Theta(d_0^2 d_1^2),$$

and $\mathbb{E}[G_1 G_1^\top]$ is permutation-invariant in the hidden index with diagonal entries

$$\mathbb{E}[(G_1 G_1^\top)_{aa}] = \frac{d_0 d_2^2}{4} + d_0 d_2 \left[(d_1 - 1) \kappa_1 + \frac{1}{2} \right] + \frac{\|\Phi\|_F^2}{4}$$

and off-diagonal entries $\mathbb{E}[(G_1 G_1^\top)_{ac}] = d_0 d_2 \kappa_2 + O(d_2 \sqrt{d_0})$ for $a \neq c$.

For G_1 the diagonal is $\Theta(d_0 d_2^2 + d_0 d_1 d_2)$ against an off-diagonal of $\Theta(d_0 d_2)$, a ratio of $\Theta(d_1)$, so $G_1 G_1^\top$ is asymptotically diagonal; for G_2 the target-free $\alpha = \Theta(d_0^2 d_1^2)$ beats the off-diagonal $\frac{d_1}{4} (\Phi \Phi^\top)_{ac}$ by $\Theta(d_0^2 d_1)$. The one nonzero off-diagonal, of order $\Theta(d_0 d_2)$, comes from a single term in the proof and traces back to ReLU's lack of odd symmetry, which is why the dominance here is asymptotic rather than exact. Simulations match every formula (Section C).

Summary and proof sketch. In all three settings the diagonal of $\mathbb{E}[GG^\top]$ outgrows the off-diagonal as the width grows (table below), so $P_t = V_t V_t^\top$ is asymptotically diagonal, $H_{\text{Muon}} = P_t^{1/2} \otimes \mathbf{I}_n$ is asymptotically block-diagonal in the Hessian's row-block structure, and (Section D) orthogonalization is asymptotically a row-wise normalization.

Orders of $\mathbb{E}[GG^\top]$ at a Gaussian initialization			
Model	Diagonal	Off-diagonal	Width
Matrix factorization	$\Theta(k^3)$	$\Theta(k^2)$	k
Deep linear, hidden layer i	$\Theta(V_i s_i)$	0	d_1, \dots, d_{L-1}
Deep linear, output layer	$\Theta(V_L s_L)$	$\Theta(\prod_{j < L} d_j)$	d_1, \dots, d_{L-1}
Two-layer ReLU, layer W_1	$\Theta(d_0 d_2^2 + d_0 d_1 d_2)$	$\Theta(d_0 d_2)$	d_1
Two-layer ReLU, layer W_2	$\Theta(d_0^2 d_1^2)$	$\Theta(d_1)$	d_1

The proofs all run the same way. We expand $\mathbb{E}[(GG^\top)_{ab}]$ by Gaussian (Wick/Isserlis) moments; a diagonal entry ($a = b$) and an off-diagonal one ($a \neq b$) differ in how many of the underlying i.i.d. Gaussian coordinates the surviving pairings tie together. On the diagonal the relevant quantities (gradient rows, columns of U , hidden neurons) share an index and are strongly correlated, so many pairings survive and their contributions accumulate over the width; off the diagonal the correlations are weak, since independent Gaussian vectors are near-orthogonal in expectation, so the cross terms cancel or contribute only lower order. This is the same effect that makes a Gaussian vector's self-inner-product grow with its dimension while its cross inner-products do not. For the deep linear network the off-diagonal cancellation is exact at hidden layers because of the sign symmetry; ReLU breaks that symmetry and leaves an $O(1)$ -per-pair residual, which still loses to the $\Theta(d_1)$ diagonal.

References

- [1] Kang An, Yuxing Liu, Rui Pan, Yi Ren, Shiqian Ma, Donald Goldfarb, and Tong Zhang. ASGO: Adaptive structured gradient optimization. *arXiv preprint arXiv:2503.20762*, 2025.
- [2] Jeremy Bernstein and Laker Newhouse. Modular duality in deep learning. *arXiv preprint arXiv:2410.21265*, 2024.
- [3] Jeremy Bernstein and Laker Newhouse. Old optimizer, new norm: An anthology. *arXiv preprint arXiv:2409.20325*, 2024.
- [4] Lizhang Chen, Jikai Li, and Qiang Liu. Muon optimizes under spectral norm constraints. *arXiv preprint arXiv:2506.15054*, 2025.
- [5] Youngmin Cho and Lawrence K. Saul. Kernel methods for deep learning. In *Advances in Neural Information Processing Systems*, volume 22, pages 342–350, 2009.
- [6] Damek Davis and Dmitriy Drusvyatskiy. When do spectral gradient updates help in deep learning? *arXiv preprint arXiv:2512.04299*, 2025.
- [7] Shenyang Deng, Zhuoli Ouyang, Tianyu Pang, Zihang Liu, Ruochen Jin, Shuhua Yu, and Yaoqing Yang. RMNP: Row-momentum normalized preconditioning for scalable matrix-based optimization. *arXiv preprint arXiv:2603.20527*, 2026.
- [8] Zhaorui Dong, Yushun Zhang, Jianfeng Yao, and Ruoyu Sun. Towards quantifying the Hessian structure of neural networks. *arXiv preprint arXiv:2505.02809*, 2025.
- [9] Chen Fan, Mark Schmidt, and Christos Thrampoulidis. Implicit bias of spectral descent and Muon on multiclass separable data. *arXiv preprint arXiv:2502.04664*, 2025.
- [10] Athanasios Glentis, Jiaxiang Li, Andi Han, and Mingyi Hong. A minimalist optimizer design for LLM pretraining. *arXiv preprint arXiv:2506.16659*, 2025.
- [11] Yufei Gu and Zeke Xie. Mano: Restriking manifold optimization for LLM training. *arXiv preprint arXiv:2601.23000*, 2026.
- [12] Leon Isserlis. On a formula for the product-moment coefficient of any order of a normal frequency distribution in any number of variables. *Biometrika*, 12(1–2):134–139, 1918.
- [13] Svante Janson. *Gaussian Hilbert Spaces*. Cambridge University Press, 1997.
- [14] Keller Jordan, Yuchen Jin, Vlado Boza, Jiacheng You, Franz Cesista, Laker Newhouse, and Jeremy Bernstein. Muon: An optimizer for hidden layers in neural networks, 2024. <https://kellerjordan.github.io/posts/muon/>.
- [15] Dmitry Kovalev. Understanding gradient orthogonalization for deep learning via non-Euclidean trust-region optimization. *arXiv preprint arXiv:2503.12645*, 2025.
- [16] Tim Tsz-Kit Lau, Qi Long, and Weijie Su. PolarGrad: A class of matrix-gradient optimizers from a unifying preconditioning perspective. *arXiv preprint arXiv:2505.21799*, 2025.

- [17] Jiayang Li and Mingyi Hong. A note on the convergence of Muon. *arXiv preprint arXiv:2502.02900*, 2025.
- [18] Chao Ma, Wenbo Gong, Meyer Scetbon, and Edward Meeds. SWAN: SGD with normalization and whitening enables stateless LLM training. In *Proceedings of the 42nd International Conference on Machine Learning (ICML)*, 2025.
- [19] Jianhao Ma, Yiding Huang, Yuejie Chi, and Yuxin Chen. Preconditioning benefits of spectral orthogonalization in Muon. *arXiv preprint arXiv:2601.13474*, 2026.
- [20] Ivan Nourdin and Giovanni Peccati. *Normal Approximations with Malliavin Calculus: From Stein’s Method to Universality*. Cambridge University Press, 2012.
- [21] Thomas Pethick, Wanyun Xie, Kimon Antonakopoulos, Zhenyu Zhu, Antonio Silveti-Falls, and Volkan Cevher. Training deep learning models with norm-constrained LMOs. In *Proceedings of the 42nd International Conference on Machine Learning (ICML)*, 2025.
- [22] Robert Price. A useful theorem for nonlinear devices having Gaussian inputs. *IRE Transactions on Information Theory*, IT-4(2):69–72, 1958.
- [23] Meyer Scetbon, Chao Ma, Wenbo Gong, and Edward Meeds. Gradient multi-normalization for stateless and scalable LLM training. *arXiv preprint arXiv:2502.06742*, 2025.
- [24] Wei Shen, Ruichuan Huang, Minhui Huang, Cong Shen, and Jiawei Zhang. On the convergence analysis of Muon. *arXiv preprint arXiv:2505.23737*, 2025.
- [25] Charles M. Stein. Estimation of the mean of a multivariate normal distribution. *The Annals of Statistics*, 9(6):1135–1151, 1981.
- [26] Weijie Su. Isotropic curvature model for understanding deep learning optimization: Is gradient orthogonalization optimal? *arXiv preprint arXiv:2511.00674*, 2025.
- [27] Bhavya Vasudeva, Puneesh Deora, Yize Zhao, Vatsal Sharan, and Christos Thrampoulidis. How Muon’s spectral design benefits generalization: A study on imbalanced data. *arXiv preprint arXiv:2510.22980*, 2025.
- [28] Shuche Wang, Fengzhuo Zhang, Jiayang Li, Cunxiao Du, Chao Du, Tianyu Pang, Zhuoran Yang, Mingyi Hong, and Vincent Y. F. Tan. Muon outperforms Adam in tail-end associative memory learning. *arXiv preprint arXiv:2509.26030*, 2025.
- [29] Ziqing Wen, Yanqi Shi, Jiahuan Wang, Ping Luo, Linbo Qiao, Dongsheng Li, and Tao Sun. SRON: State-free LLM training via row-wise gradient normalization. In *OpenReview preprint*, 2025. <https://openreview.net/forum?id=BtQLBWr6zI>.
- [30] Gian-Carlo Wick. The evaluation of the collision matrix. *Physical Review*, 80(2):268–272, 1950.
- [31] Ruihan Xu, Jiajin Li, and Yiping Lu. On the width scaling of neural optimizers under matrix operator norms I: Row/column normalization and hyperparameter transfer. *arXiv preprint arXiv:2603.09952*, 2026.

- [32] Thomas T. Zhang, Behrad Moniri, Ansh Nagwekar, Faraz Rahman, Anton Xue, Hamed Hasani, and Nikolai Matni. On the concurrence of layer-wise preconditioning methods and provable feature learning. *arXiv preprint arXiv:2502.01763*, 2025.
- [33] Yushun Zhang, Congliang Chen, Tian Ding, Ziniu Li, Ruoyu Sun, and Zhi-Quan Luo. Why transformers need Adam: A Hessian perspective. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- [34] Yushun Zhang, Congliang Chen, Ziniu Li, Tian Ding, Chenwei Wu, Yinyu Ye, Zhi-Quan Luo, and Ruoyu Sun. Adam-mini: Use fewer learning rates to gain more. *arXiv preprint arXiv:2406.16793*, 2024.

Contents

1	Introduction	1
2	Related Work	2
3	Preliminaries	3
4	Main Results	4
A	Symmetric Matrix Factorization	9
A.1	Setup	9
A.2	Standing assumption	9
A.3	Main results restated	9
A.4	Preliminary lemmas	10
A.5	Proofs	10
B	Deep Linear Networks	12
B.1	Setup	12
B.2	Standing assumption	12
B.3	Main results restated	12
B.4	Preliminary lemmas	13
B.5	Proofs	14
C	Two-Layer ReLU Networks	15
C.1	Setup	15
C.2	Standing assumption	16
C.3	Main results restated	16
C.4	Preliminary lemmas	16
C.5	Proof of the gradient self outer-product of W_2	17
C.6	Proof of the gradient self outer-product of W_1	18
C.7	Numerical verification	20
C.8	Discussion	21

D Orthogonalization Is Asymptotically a Row Normalization

22

Appendix A. Symmetric Matrix Factorization

A.1. Setup

Let d, k, r be positive integers with $k \geq r$. The over-parameterized symmetric matrix-factorization loss is

$$\mathcal{L}(\mathbf{U}) := \frac{1}{4} \left\| \mathbf{U}\mathbf{U}^\top - \mathbf{M}^* \right\|_{\text{F}}^2, \quad \mathbf{U} \in \mathbb{R}^{d \times k}, \quad (5)$$

with target $\mathbf{M}^* \in \mathbb{R}^{d \times d}$ symmetric positive semidefinite of rank r . A direct computation gives the gradient

$$\mathbf{G} := \nabla \mathcal{L}(\mathbf{U}) = (\mathbf{U}\mathbf{U}^\top - \mathbf{M}^*)\mathbf{U} \in \mathbb{R}^{d \times k}, \quad (6)$$

and the object of interest is the $d \times d$ matrix $\mathbb{E}[\mathbf{G}\mathbf{G}^\top]$.

A.2. Standing assumption

Assumption 2 *The entries U_{ij} are i.i.d. $\mathcal{N}(0, 1)$, and \mathbf{M}^* is deterministic and independent of \mathbf{U} . Hence $\mathbf{W} := \mathbf{U}\mathbf{U}^\top$ is Wishart with k degrees of freedom and identity scale, $\mathbf{W} \sim \mathcal{W}_d(k, \mathbf{I}_d)$.*

A.3. Main results restated

Theorem (Theorem 1, restated). *Under Assumption 2,*

$$\mathbb{E}[\mathbf{G}\mathbf{G}^\top] = \alpha_3(d, k) \mathbf{I}_d - 2\alpha_2(d, k) \mathbf{M}^* + k(\mathbf{M}^*)^2, \quad (7)$$

where the scalars are the diagonal entries of the first three Wishart moments,

$$\alpha_1(d, k) = k, \quad (8)$$

$$\alpha_2(d, k) = k(k + d + 1), \quad (9)$$

$$\alpha_3(d, k) = k^3 + 3(d + 1)k^2 + (d^2 + 3d + 4)k. \quad (10)$$

Corollary 1 *Under Assumption 2, $\mathbb{E}[\mathbf{G}\mathbf{G}^\top]$ is invariant under any orthogonal change of basis fixing \mathbf{M}^* , with*

$$\mathbb{E}[(\mathbf{G}\mathbf{G}^\top)_{ii}] = \alpha_3 - 2\alpha_2 \mathbf{M}_{ii}^* + k(\mathbf{M}_{ii}^*)^2, \quad (11)$$

$$\mathbb{E}[(\mathbf{G}\mathbf{G}^\top)_{ij}] = -2\alpha_2 \mathbf{M}_{ij}^* + k(\mathbf{M}_{ij}^*)^2 \quad (i \neq j). \quad (12)$$

For fixed d , \mathbf{M}^* as $k \rightarrow \infty$, $\mathbb{E}[(\mathbf{G}\mathbf{G}^\top)_{ii}] = \Theta(k^3)$ and $\mathbb{E}[(\mathbf{G}\mathbf{G}^\top)_{ij}] = \Theta(k^2)$ whenever $\mathbf{M}_{ij}^* \neq 0$, so the dominance ratio is $\Theta(k/\mathbf{M}_{ij}^*)$.

A.4. Preliminary lemmas

Lemma 1 *Let $\mathbf{X} \in \mathbb{R}^{m \times n}$ have i.i.d. $\mathcal{N}(0, 1)$ entries. Then for deterministic $\mathbf{M} \in \mathbb{R}^{n \times n}$ and $\mathbf{N} \in \mathbb{R}^{m \times m}$,*

$$\mathbb{E}[\mathbf{X}\mathbf{M}\mathbf{X}^\top] = \text{tr}(\mathbf{M}) \mathbf{I}_m, \quad \mathbb{E}[\mathbf{X}^\top \mathbf{N} \mathbf{X}] = \text{tr}(\mathbf{N}) \mathbf{I}_n. \quad (13)$$

Proof. Using $\mathbb{E}[\mathbf{X}_{ap}\mathbf{X}_{bq}] = \delta_{ab}\delta_{pq}$,

$$\mathbb{E}[(\mathbf{X}\mathbf{M}\mathbf{X}^\top)_{ab}] = \sum_{p,q} M_{pq} \delta_{ab}\delta_{pq} = \delta_{ab} \text{tr}(\mathbf{M}),$$

which is the first identity; the second follows by applying it to \mathbf{X}^\top (entries i.i.d. $\mathcal{N}(0, 1)$). ■

Lemma 2 *Let $\mathbf{W} = \mathbf{U}\mathbf{U}^\top$ with \mathbf{U} as in Assumption 2. For every $m \in \mathbb{N}$ there is a scalar $\alpha_m(d, k)$ with*

$$\mathbb{E}[\mathbf{W}^m] = \alpha_m(d, k) \mathbf{I}_d, \quad \alpha_m(d, k) = \frac{1}{d} \mathbb{E}[\text{tr}(\mathbf{W}^m)]. \quad (14)$$

Proof. For $\mathbf{Q} \in O(d)$, $\mathbf{Q}\mathbf{U} \stackrel{d}{=} \mathbf{U}$, hence $\mathbf{Q}\mathbf{W}^m\mathbf{Q}^\top \stackrel{d}{=} \mathbf{W}^m$ and

$$\mathbb{E}[\mathbf{W}^m] = \mathbf{Q} \mathbb{E}[\mathbf{W}^m] \mathbf{Q}^\top \quad \text{for all } \mathbf{Q} \in O(d).$$

By Schur's lemma $\mathbb{E}[\mathbf{W}^m] = \alpha_m \mathbf{I}_d$; taking the trace gives $d\alpha_m = \mathbb{E}[\text{tr}(\mathbf{W}^m)]$. ■

Lemma 3 (Wick/Isserlis) *Let X_1, \dots, X_{2n} be jointly centered Gaussian. Then*

$$\mathbb{E}[X_1 \cdots X_{2n}] = \sum_P \prod_{\{i,j\} \in P} \mathbb{E}[X_i X_j], \quad (15)$$

the sum over all pairings P of $\{1, \dots, 2n\}$; odd moments vanish [12, 30].

Proof. Differentiate $\mathbb{E}[\exp(\mathbf{t}^\top \mathbf{X})] = \exp(\frac{1}{2} \mathbf{t}^\top \Sigma \mathbf{t})$ in t_1, \dots, t_{2n} at $\mathbf{t} = \mathbf{0}$; cf. Janson [13, §1.4]. ■

A.5. Proofs

Lemma 4 *With $\mathbf{W} := \mathbf{U}\mathbf{U}^\top$ and $\mathbf{B} := \mathbf{W} - \mathbf{M}^*$,*

$$\mathbf{G}\mathbf{G}^\top = \mathbf{B}\mathbf{W}\mathbf{B} = \mathbf{W}^3 - \mathbf{W}^2\mathbf{M}^* - \mathbf{M}^*\mathbf{W}^2 + \mathbf{M}^*\mathbf{W}\mathbf{M}^*. \quad (16)$$

Proof. Since $\mathbf{B} = \mathbf{B}^\top$,

$$\mathbf{G}\mathbf{G}^\top = \mathbf{B}\mathbf{U}(\mathbf{B}\mathbf{U})^\top = \mathbf{B}\mathbf{U}\mathbf{U}^\top\mathbf{B}^\top = \mathbf{B}\mathbf{W}\mathbf{B},$$

and substituting $\mathbf{B} = \mathbf{W} - \mathbf{M}^*$ and expanding gives (16). ■

Lemma 5 *Under Assumption 2, $\mathbb{E}[\mathbf{W}] = k \mathbf{I}_d$, hence $\alpha_1(d, k) = k$.*

Proof. Let \mathbf{u}_a denote the a -th column of \mathbf{U} . Then

$$\mathbb{E}[\mathbf{W}] = \mathbb{E}\left[\sum_{a=1}^k \mathbf{u}_a \mathbf{u}_a^\top\right] = \sum_{a=1}^k \mathbb{E}[\mathbf{u}_a \mathbf{u}_a^\top] = k \mathbf{I}_d. \quad \blacksquare$$

Lemma 6 Under Assumption 2, $\mathbb{E}[\mathbf{W}^2] = k(k+d+1)\mathbf{I}_d$, hence $\alpha_2(d, k) = k(k+d+1)$.

Proof. By Lemma 2, $\mathbb{E}[\mathbf{W}^2] = \alpha_2\mathbf{I}_d$. With columns $\mathbf{u}_a \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$,

$$\mathbf{W}^2 = \sum_{a,b=1}^k \langle \mathbf{u}_a, \mathbf{u}_b \rangle \mathbf{u}_a \mathbf{u}_b^\top.$$

Split into the k terms with $a = b$ and the $k(k-1)$ with $a \neq b$. By Lemma 3 and independence,

$$\begin{aligned} \mathbb{E}[\|\mathbf{u}_a\|^2 \mathbf{u}_a \mathbf{u}_a^\top] &= (d+2)\mathbf{I}_d, \\ \mathbb{E}[\langle \mathbf{u}_a, \mathbf{u}_b \rangle \mathbf{u}_a \mathbf{u}_b^\top] &= \sum_p \mathbb{E}[u_{a,p} \mathbf{u}_a] \mathbb{E}[u_{b,p} \mathbf{u}_b^\top] = \sum_p \mathbf{e}_p \mathbf{e}_p^\top = \mathbf{I}_d. \end{aligned}$$

Summing, $\mathbb{E}[\mathbf{W}^2] = k(d+2)\mathbf{I}_d + k(k-1)\mathbf{I}_d = k(k+d+1)\mathbf{I}_d$. \blacksquare

Lemma 7 Under Assumption 2, $\mathbb{E}[\mathbf{W}^3] = [k^3 + 3(d+1)k^2 + (d^2 + 3d + 4)k]\mathbf{I}_d$, hence $\alpha_3(d, k)$ as in (10).

Proof. By Lemma 2, $\alpha_3 = \frac{1}{d}\mathbb{E}[\text{tr}(\mathbf{W}^3)]$, and with $\mathbf{W} = \sum_a \mathbf{u}_a \mathbf{u}_a^\top$ and cyclicity of the trace,

$$\text{tr}(\mathbf{W}^3) = \sum_{a,b,c} \langle \mathbf{u}_a, \mathbf{u}_b \rangle \langle \mathbf{u}_b, \mathbf{u}_c \rangle \langle \mathbf{u}_c, \mathbf{u}_a \rangle.$$

Partition $(a, b, c) \in \{1, \dots, k\}^3$ by multiplicity pattern:

$$\begin{aligned} \#\{a = b = c\} = k &: & \mathbb{E}[\|\mathbf{u}_a\|^6] &= d(d+2)(d+4), \\ \#\{|\{a, b, c\}| = 2\} = 3k(k-1) &: & \mathbb{E}[\|\mathbf{u}_a\|^2 \langle \mathbf{u}_a, \mathbf{u}_c \rangle^2] &= \mathbb{E}[\|\mathbf{u}_a\|^4] = d(d+2), \\ \#\{a, b, c \text{ distinct}\} = k(k-1)(k-2) &: & \sum_{p,q,r} \delta_{pr} \delta_{pq} \delta_{qr} &= d, \end{aligned}$$

where for the middle pattern, e.g. $a = b \neq c$ gives $\|\mathbf{u}_a\|^2 \langle \mathbf{u}_a, \mathbf{u}_c \rangle^2$ and, conditioning on \mathbf{u}_a , $\langle \mathbf{u}_a, \mathbf{u}_c \rangle \mid \mathbf{u}_a \sim \mathcal{N}(0, \|\mathbf{u}_a\|^2)$; the other two cases are symmetric. Hence

$$d\alpha_3 = k d(d+2)(d+4) + 3k(k-1) d(d+2) + k(k-1)(k-2) d,$$

$$\alpha_3 = k(d+2)(d+4) + 3k(k-1)(d+2) + k(k-1)(k-2) = k^3 + 3(d+1)k^2 + (d^2 + 3d + 4)k. \quad \blacksquare$$

Proof of Theorem 1. Take expectations in (16). Since \mathbf{M}^* is deterministic, Lemma 2 gives

$$\begin{aligned} \mathbb{E}[\mathbf{W}^3] &= \alpha_3\mathbf{I}_d, \\ \mathbb{E}[\mathbf{W}^2\mathbf{M}^*] &= \mathbb{E}[\mathbf{M}^*\mathbf{W}^2] = \alpha_2\mathbf{M}^*, \\ \mathbb{E}[\mathbf{M}^*\mathbf{W}\mathbf{M}^*] &= \mathbf{M}^* \mathbb{E}[\mathbf{W}] \mathbf{M}^* = k(\mathbf{M}^*)^2, \end{aligned}$$

so summing gives (7); the values of $\alpha_1, \alpha_2, \alpha_3$ are Lemmas 5 to 7. Corollary 1 follows by reading off entries: with d, \mathbf{M}^* fixed, $\alpha_3 = \Theta(k^3)$, $\alpha_2 = \Theta(k^2)$, $k(\mathbf{M}^*)^2 = \Theta(k)$. \blacksquare

Remark 1 This computation answers the question raised by Ma et al. [19] of how gradient orthogonalization meets the curvature of training: the $\Theta(k)$ gap between the diagonal and off-diagonal scales is the regime in which orthogonalization, which suppresses the isotropic $\Theta(k^3)\mathbf{I}_d$ part, acts as a real preconditioner relative to the target-aligned signal \mathbf{M}^* .

Appendix B. Deep Linear Networks

B.1. Setup

Let $L \geq 1$ and d_0, \dots, d_L be positive integers. The depth- L linear network

$$f(\mathbf{x}) = \mathbf{W}_L \mathbf{W}_{L-1} \cdots \mathbf{W}_1 \mathbf{x}, \quad \mathbf{W}_i \in \mathbb{R}^{d_i \times d_{i-1}}, \quad (17)$$

learns a deterministic linear target $\Phi \in \mathbb{R}^{d_L \times d_0}$ under the population squared loss

$$\mathcal{L}(\mathbf{W}_1, \dots, \mathbf{W}_L) := \frac{1}{2} \mathbb{E}_{\mathbf{x}} \|\mathbf{W}_L \cdots \mathbf{W}_1 \mathbf{x} - \Phi \mathbf{x}\|^2. \quad (18)$$

Define the layer- i forward, backward and residual factors

$$\mathbf{B}_i := \mathbf{W}_{i-1} \cdots \mathbf{W}_1 \in \mathbb{R}^{d_{i-1} \times d_0}, \quad \mathbf{A}_i := \mathbf{W}_L \cdots \mathbf{W}_{i+1} \in \mathbb{R}^{d_L \times d_i}, \quad \mathbf{R} := \mathbf{W}_L \cdots \mathbf{W}_1 - \Phi, \quad (19)$$

with $\mathbf{B}_1 := \mathbf{I}_{d_0}$ and $\mathbf{A}_L := \mathbf{I}_{d_L}$. The layer- i gradient is

$$\mathbf{G}_i := \frac{\partial \mathcal{L}}{\partial \mathbf{W}_i} = \mathbf{A}_i^\top \mathbf{R} \mathbf{B}_i^\top \in \mathbb{R}^{d_i \times d_{i-1}}. \quad (20)$$

The closed form below uses two pairs of dimension-only scalar sequences. For the backward factors $\{\mathbf{B}_i\}$,

$$r_1 = 1, \quad v_1 = d_0^2, \quad (21)$$

$$r_i = v_{i-1} + (d_{i-1} + 1)d_{i-2}r_{i-1}, \quad v_i = d_{i-1}^2 v_{i-1} + 2d_{i-1}d_{i-2}r_{i-1} \quad (i \geq 2); \quad (22)$$

for the forward factors $\{\mathbf{A}_i\}$,

$$s_L = 1, \quad u_L = d_L^2, \quad (23)$$

$$s_{i-1} = u_i + (d_{i-1} + 1)d_i s_i, \quad u_{i-1} = d_{i-1}^2 u_i + 2d_{i-1}d_i s_i \quad (i \leq L); \quad (24)$$

and the composite Frobenius coefficient is

$$V_i := d_{i-1} r_i = \mathbb{E} \left[\left\| \mathbf{B}_i \mathbf{B}_i^\top \right\|_{\text{F}}^2 \right]. \quad (25)$$

B.2. Standing assumption

Assumption 3 *The entries $(\mathbf{W}_i)_{ab}$ are i.i.d. $\mathcal{N}(0, 1)$ across all i, a, b ; $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_0})$; and Φ is deterministic and independent of $\mathbf{W}_1, \dots, \mathbf{W}_L$. Expectations $\mathbb{E}[\cdot]$ are over $\mathbf{W}_1, \dots, \mathbf{W}_L$ (\mathbf{x} has been integrated out in \mathcal{L}).*

B.3. Main results restated

Theorem (Theorem 2, restated). *Under Assumption 3, for every layer $i \in \{1, \dots, L\}$,*

$$\mathbb{E}[\mathbf{G}_i \mathbf{G}_i^\top] = V_i s_i \mathbf{I}_{d_i} + \mathbf{T}_i, \quad \mathbf{T}_i = \begin{cases} \left(\prod_{j=1}^{L-1} d_j \right) \Phi \Phi^\top, & i = L, \\ \left(\prod_{j=1, j \neq i}^{L-1} d_j \right) \|\Phi\|_{\text{F}}^2 \mathbf{I}_{d_i}, & 1 \leq i < L, \end{cases} \quad (26)$$

with V_i, s_i from (21)–(25).

Corollary 2 Under Assumption 3, the diagonal entries of $\mathbb{E}[\mathbf{G}_i \mathbf{G}_i^\top]$ are

$$\mathbb{E}[(\mathbf{G}_i \mathbf{G}_i^\top)_{aa}] = \begin{cases} V_i s_i + (\prod_{j \neq i, 1 \leq j \leq L-1} d_j) \|\Phi\|_F^2, & 1 \leq i < L, \\ V_L s_L + (\prod_{j=1}^{L-1} d_j) (\Phi \Phi^\top)_{aa}, & i = L, \end{cases} \quad (27)$$

and, for $a \neq b$, the off-diagonal entries are $\mathbb{E}[(\mathbf{G}_i \mathbf{G}_i^\top)_{ab}] = 0$ if $1 \leq i < L$ and $(\prod_{j=1}^{L-1} d_j) (\Phi \Phi^\top)_{ab}$ if $i = L$. In particular $\mathbb{E}[\mathbf{G}_i \mathbf{G}_i^\top]$ is a scalar multiple of \mathbf{I}_{d_i} for every $i < L$, and unrolling (21)–(24) gives $V_i s_i \asymp d_{i-1} (\prod_{j=0}^{i-2} d_j)^2 (\prod_{j=i+1}^L d_j)^2$.

B.4. Preliminary lemmas

Lemma 8 Same as Lemma 1.

Lemma 9 (Wick/Isserlis) Same as Lemma 3.

Lemma 10 Let $\mathbf{X} \in \mathbb{R}^{m \times n}$ have i.i.d. $\mathcal{N}(0, 1)$ entries and let $\mathbf{M} \in \mathbb{R}^{n \times n}$, $\mathbf{N} \in \mathbb{R}^{m \times m}$ be deterministic symmetric. Then

$$\mathbb{E}[\mathbf{X} \mathbf{M} \mathbf{X}^\top \mathbf{X} \mathbf{M} \mathbf{X}^\top] = (\text{tr}(\mathbf{M})^2 + (m+1) \|\mathbf{M}\|_F^2) \mathbf{I}_m, \quad (28)$$

$$\mathbb{E}[\mathbf{X}^\top \mathbf{N} \mathbf{X} \mathbf{X}^\top \mathbf{N} \mathbf{X}] = (\text{tr}(\mathbf{N})^2 + (n+1) \|\mathbf{N}\|_F^2) \mathbf{I}_n. \quad (29)$$

Proof. We prove (28); (29) follows by applying it to \mathbf{X}^\top . For $a, b, c, d \in \{1, \dots, m\}$, Lemma 9 gives

$$\mathbb{E}[\mathbf{X}_{ap} \mathbf{X}_{bq} \mathbf{X}_{cr} \mathbf{X}_{ds}] = \delta_{ab} \delta_{pq} \delta_{cd} \delta_{rs} + \delta_{ac} \delta_{pr} \delta_{bd} \delta_{qs} + \delta_{ad} \delta_{ps} \delta_{bc} \delta_{qr},$$

and contracting against $\mathbf{M}_{pq} \mathbf{M}_{rs}$ with $\mathbf{M} = \mathbf{M}^\top$,

$$\mathbb{E}[(\mathbf{X} \mathbf{M} \mathbf{X}^\top)_{ab} (\mathbf{X} \mathbf{M} \mathbf{X}^\top)_{cd}] = \delta_{ab} \delta_{cd} \text{tr}(\mathbf{M})^2 + (\delta_{ac} \delta_{bd} + \delta_{ad} \delta_{bc}) \|\mathbf{M}\|_F^2.$$

Summing over $b = c$,

$$\mathbb{E}[(\mathbf{X} \mathbf{M} \mathbf{X}^\top \mathbf{X} \mathbf{M} \mathbf{X}^\top)_{ac}] = \sum_b \mathbb{E}[(\mathbf{X} \mathbf{M} \mathbf{X}^\top)_{ab} (\mathbf{X} \mathbf{M} \mathbf{X}^\top)_{bc}] = \delta_{ac} (\text{tr}(\mathbf{M})^2 + (m+1) \|\mathbf{M}\|_F^2). \quad \blacksquare$$

Lemma 11 Let $\mathbf{X} \in \mathbb{R}^{m \times n}$ have i.i.d. $\mathcal{N}(0, 1)$ entries and $\mathbf{M} \in \mathbb{R}^{n \times n}$ be deterministic symmetric. Then

$$\mathbb{E}[\text{tr}(\mathbf{X} \mathbf{M} \mathbf{X}^\top)^2] = m^2 \text{tr}(\mathbf{M})^2 + 2m \|\mathbf{M}\|_F^2. \quad (30)$$

Proof. Expanding,

$$\text{tr}(\mathbf{X} \mathbf{M} \mathbf{X}^\top)^2 = \sum_{a, a', p, q, p', q'} \mathbf{X}_{ap} \mathbf{M}_{pq} \mathbf{X}_{aq} \mathbf{X}_{a'p'} \mathbf{M}_{p'q'} \mathbf{X}_{a'q'};$$

by Lemma 9 the three pair contractions of $\mathbb{E}[\mathbf{X}_{ap} \mathbf{X}_{aq} \mathbf{X}_{a'p'} \mathbf{X}_{a'q'}]$ contribute (after summing over $\mathbf{M} \mathbf{M}$) $m^2 \text{tr}(\mathbf{M})^2$, $m \|\mathbf{M}\|_F^2$ and $m \text{tr}(\mathbf{M}^2) = m \|\mathbf{M}\|_F^2$, which sum to (30). \blacksquare

B.5. Proofs

The proof of Theorem 2 has three steps: a \mathbf{W}_i -expansion; recursive moments of $\mathbf{B}_i\mathbf{B}_i^\top$ and $\mathbf{A}_i^\top\mathbf{A}_i$; and the target-dependent term.

Lemma 12 *Under Assumption 3, for every i ,*

$$\mathbb{E}[\mathbf{G}_i\mathbf{G}_i^\top \mid \mathbf{A}_i, \mathbf{B}_i] = \text{tr}((\mathbf{B}_i\mathbf{B}_i^\top)^2) (\mathbf{A}_i^\top\mathbf{A}_i)^2 + \mathbf{A}_i^\top\Phi\mathbf{B}_i^\top\mathbf{B}_i\Phi^\top\mathbf{A}_i. \quad (31)$$

Proof. Substituting $\mathbf{R} = \mathbf{A}_i\mathbf{W}_i\mathbf{B}_i - \Phi$ into (20),

$$\mathbf{G}_i\mathbf{G}_i^\top = (\mathbf{A}_i^\top\mathbf{A}_i)\mathbf{W}_i(\mathbf{B}_i\mathbf{B}_i^\top\mathbf{B}_i\mathbf{B}_i^\top)\mathbf{W}_i^\top(\mathbf{A}_i^\top\mathbf{A}_i) + (\text{cross}) + \mathbf{A}_i^\top\Phi\mathbf{B}_i^\top\mathbf{B}_i\Phi^\top\mathbf{A}_i,$$

where (cross) collects the two terms linear in \mathbf{W}_i ; since $\mathbf{W}_i \stackrel{d}{=} -\mathbf{W}_i$ is independent of $(\mathbf{A}_i, \mathbf{B}_i)$, $\mathbb{E}[(\text{cross}) \mid \mathbf{A}_i, \mathbf{B}_i] = \mathbf{0}$. By Lemma 8 with $\mathbf{K} = \mathbf{B}_i\mathbf{B}_i^\top\mathbf{B}_i\mathbf{B}_i^\top$,

$$\mathbb{E}_{\mathbf{W}_i}[\mathbf{W}_i\mathbf{K}\mathbf{W}_i^\top \mid \mathbf{A}_i, \mathbf{B}_i] = \text{tr}(\mathbf{K})\mathbf{I}_{d_i} = \text{tr}((\mathbf{B}_i\mathbf{B}_i^\top)^2)\mathbf{I}_{d_i},$$

which yields (31). ■

Lemma 13 *Under Assumption 3, for every $i \in \{1, \dots, L\}$,*

$$\mathbb{E}[(\mathbf{B}_i\mathbf{B}_i^\top)^2] = r_i\mathbf{I}_{d_{i-1}}, \quad \mathbb{E}[\text{tr}(\mathbf{B}_i\mathbf{B}_i^\top)^2] = v_i, \quad (32)$$

with (r_i, v_i) from (21)–(22).

Proof. At $i = 1$, $\mathbf{B}_1 = \mathbf{I}_{d_0}$, so $(\mathbf{B}_1\mathbf{B}_1^\top)^2 = \mathbf{I}_{d_0}$ and $\text{tr}(\mathbf{B}_1\mathbf{B}_1^\top)^2 = d_0^2$, matching $r_1 = 1$, $v_1 = d_0^2$. For $i \geq 2$, $\mathbf{B}_i = \mathbf{W}_{i-1}\mathbf{B}_{i-1}$ with $\mathbf{W}_{i-1} \in \mathbb{R}^{d_{i-1} \times d_{i-2}}$ independent of \mathbf{B}_{i-1} , so $\mathbf{B}_i\mathbf{B}_i^\top = \mathbf{W}_{i-1}\mathbf{M}\mathbf{W}_{i-1}^\top$ with $\mathbf{M} := \mathbf{B}_{i-1}\mathbf{B}_{i-1}^\top$. By Lemma 10,

$$\mathbb{E}[(\mathbf{B}_i\mathbf{B}_i^\top)^2 \mid \mathbf{B}_{i-1}] = (\text{tr}(\mathbf{M})^2 + (d_{i-1} + 1)\|\mathbf{M}\|_{\mathbb{F}}^2)\mathbf{I}_{d_{i-1}},$$

and taking \mathbb{E} over \mathbf{B}_{i-1} with $\mathbb{E}[\text{tr}(\mathbf{M}^2)] = d_{i-2}r_{i-1}$ gives $r_i = v_{i-1} + (d_{i-1} + 1)d_{i-2}r_{i-1}$. Likewise, by Lemma 11,

$$\mathbb{E}[\text{tr}(\mathbf{B}_i\mathbf{B}_i^\top)^2 \mid \mathbf{B}_{i-1}] = d_{i-1}^2\text{tr}(\mathbf{M})^2 + 2d_{i-1}\|\mathbf{M}\|_{\mathbb{F}}^2,$$

so $v_i = d_{i-1}^2v_{i-1} + 2d_{i-1}d_{i-2}r_{i-1}$. ■

Lemma 14 *Under Assumption 3, for every $i \in \{0, \dots, L\}$,*

$$\mathbb{E}[(\mathbf{A}_i^\top\mathbf{A}_i)^2] = s_i\mathbf{I}_{d_i}, \quad \mathbb{E}[\text{tr}(\mathbf{A}_i^\top\mathbf{A}_i)^2] = u_i, \quad (33)$$

with (s_i, u_i) from (23)–(24).

Proof. At $i = L$, $\mathbf{A}_L = \mathbf{I}_{d_L}$, matching $s_L = 1$, $u_L = d_L^2$. For $i < L$, $\mathbf{A}_i = \mathbf{A}_{i+1}\mathbf{W}_{i+1}$ with $\mathbf{W}_{i+1} \in \mathbb{R}^{d_{i+1} \times d_i}$ independent of \mathbf{A}_{i+1} , so $\mathbf{A}_i^\top\mathbf{A}_i = \mathbf{W}_{i+1}^\top\mathbf{N}\mathbf{W}_{i+1}$ with $\mathbf{N} := \mathbf{A}_{i+1}^\top\mathbf{A}_{i+1}$. By Lemma 10 (right identity),

$$\mathbb{E}[(\mathbf{A}_i^\top\mathbf{A}_i)^2 \mid \mathbf{A}_{i+1}] = (\text{tr}(\mathbf{N})^2 + (d_i + 1)\|\mathbf{N}\|_{\mathbb{F}}^2)\mathbf{I}_{d_i},$$

so $s_i = u_{i+1} + (d_i + 1)d_{i+1}s_{i+1}$; the recursion for u_i is that of v_i with $d_{i-1} \leftrightarrow d_i$ and $d_{i-2} \leftrightarrow d_{i+1}$, giving $u_{i-1} = d_{i-1}^2u_i + 2d_{i-1}d_i s_i$. ■

Lemma 15 Under Assumption 3, $\mathbb{E}[\mathbf{A}_i^\top \Phi \mathbf{B}_i^\top \mathbf{B}_i \Phi^\top \mathbf{A}_i] = \mathbf{T}_i$, with \mathbf{T}_i as in (26).

Proof. By independence of $\mathbf{A}_i, \mathbf{B}_i$ and Lemma 8 iterated through $\mathbf{B}_i = \mathbf{W}_{i-1} \cdots \mathbf{W}_1$,

$$\mathbb{E}[\mathbf{B}_i^\top \mathbf{B}_i] = \left(\prod_{j=1}^{i-1} d_j \right) \mathbf{I}_{d_0} \quad \implies \quad \mathbb{E}[\mathbf{A}_i^\top \Phi \mathbf{B}_i^\top \mathbf{B}_i \Phi^\top \mathbf{A}_i] = \left(\prod_{j=1}^{i-1} d_j \right) \mathbb{E}[\mathbf{A}_i^\top \Phi \Phi^\top \mathbf{A}_i]$$

(the empty product is 1 when $i = 1$). Let $h_i(\mathbf{M}) := \mathbb{E}[\mathbf{A}_i^\top \mathbf{M} \mathbf{A}_i]$ for symmetric $\mathbf{M} \in \mathbb{R}^{d_L \times d_L}$. Then $h_L(\mathbf{M}) = \mathbf{M}$, and for $i < L$, by Lemma 8,

$$h_i(\mathbf{M}) = \mathbb{E}[\text{tr}(\mathbf{A}_{i+1}^\top \mathbf{M} \mathbf{A}_{i+1})] \mathbf{I}_{d_i} = \text{tr}(h_{i+1}(\mathbf{M})) \mathbf{I}_{d_i}.$$

Iterating,

$$h_i(\mathbf{M}) = \begin{cases} \mathbf{M}, & i = L, \\ \left(\prod_{j=i+1}^{L-1} d_j \right) \text{tr}(\mathbf{M}) \mathbf{I}_{d_i}, & i < L, \end{cases}$$

with $\prod_{j=L}^{L-1} d_j = 1$. Taking $\mathbf{M} = \Phi \Phi^\top$ (so $\text{tr}(\mathbf{M}) = \|\Phi\|_F^2$) and multiplying by $\prod_{j=1}^{i-1} d_j$ gives \mathbf{T}_i . ■

Proof of Theorem 2. By Lemma 12,

$$\mathbb{E}[\mathbf{G}_i \mathbf{G}_i^\top] = \mathbb{E}[\text{tr}((\mathbf{B}_i \mathbf{B}_i^\top)^2) (\mathbf{A}_i^\top \mathbf{A}_i)^2] + \mathbb{E}[\mathbf{A}_i^\top \Phi \mathbf{B}_i^\top \mathbf{B}_i \Phi^\top \mathbf{A}_i].$$

Since $\mathbf{A}_i, \mathbf{B}_i$ depend on disjoint sets of \mathbf{W}_j , they are independent, so by Lemmas 13 and 14,

$$\mathbb{E}[\text{tr}((\mathbf{B}_i \mathbf{B}_i^\top)^2) (\mathbf{A}_i^\top \mathbf{A}_i)^2] = \mathbb{E}[\text{tr}((\mathbf{B}_i \mathbf{B}_i^\top)^2)] \mathbb{E}[(\mathbf{A}_i^\top \mathbf{A}_i)^2] = d_{i-1} r_i \cdot s_i \mathbf{I}_{d_i} = V_i s_i \mathbf{I}_{d_i},$$

and the second term equals \mathbf{T}_i by Lemma 15. Corollary 2 follows by reading off entries of (26): for $i < L$, \mathbf{T}_i is a multiple of \mathbf{I}_{d_i} so the off-diagonal vanishes; for $i = L$, \mathbf{T}_L contributes $(\Phi \Phi^\top)_{ab}$ entry-wise. The leading order of $V_i s_i$ follows by unrolling (21)–(24). ■

Appendix C. Two-Layer ReLU Networks

C.1. Setup

Let d_0, d_1, d_2 be positive integers. The two-layer ReLU network

$$f(x) = W_2 \sigma(W_1 x), \quad W_1 \in \mathbb{R}^{d_1 \times d_0}, \quad W_2 \in \mathbb{R}^{d_2 \times d_1}, \quad (34)$$

acts on $x \in \mathbb{R}^{d_0}$ via $\sigma(z) = \max(z, 0)$, $\sigma'(z) = \mathbf{1}[z > 0]$. Write $h(x) := \sigma(W_1 x) \in \mathbb{R}^{d_1}$ and $w_a \in \mathbb{R}^{d_0}$ for the a -th row of W_1 . The population loss against a deterministic linear target $\Phi \in \mathbb{R}^{d_2 \times d_0}$ is

$$\mathcal{L}(W_1, W_2) := \frac{1}{2} \mathbb{E}_x \|W_2 \sigma(W_1 x) - \Phi x\|^2, \quad (35)$$

and $G_i := \partial \mathcal{L} / \partial W_i$, $i \in \{1, 2\}$. For independent $u, v \sim \mathcal{N}(0, \mathbf{I}_{d_0})$ let $\theta := \angle(u, v) = \arccos(\langle u, v \rangle / (\|u\| \|v\|)) \in [0, \pi]$, and define

$$\kappa_1 := \mathbb{E} \left[\frac{\sin^2 \theta + (\pi - \theta)^2}{4\pi^2} + \frac{\sin \theta \cos \theta (\pi - \theta)}{2\pi^2} \right], \quad (36)$$

$$\kappa_2 := \frac{\mathbb{E}[\sin \theta (\pi - \theta)]}{2\pi^2}. \quad (37)$$

C.2. Standing assumption

Assumption 4 $(W_1)_{ij}$ and $(W_2)_{ij}$ are i.i.d. $\mathcal{N}(0, 1)$; $x \sim \mathcal{N}(0, \mathbf{I}_{d_0})$; Φ is deterministic and independent of W_1, W_2 . Expectations $\mathbb{E}[\cdot]$ are over W_1, W_2 (x has been integrated out in \mathcal{L}).

C.3. Main results restated

Theorem (Theorem 3, restated). Under Assumption 4, with $H := \mathbb{E}_x[h(x)h(x)^\top]$ and $\alpha := \mathbb{E}_{W_1}[\text{tr}(H^2)]$,

$$\mathbb{E}[G_2 G_2^\top] = \alpha \mathbf{I}_{d_2} + \frac{d_1}{4} \Phi \Phi^\top, \quad \alpha = \Theta(d_0^2 d_1^2) \text{ as } d_0, d_1 \rightarrow \infty, \quad (38)$$

and $\mathbb{E}[G_1 G_1^\top]$ is permutation-invariant in the hidden index with, as $d_0, d_2 \rightarrow \infty$,

$$\mathbb{E}[(G_1 G_1^\top)_{aa}] = \frac{d_0 d_2^2}{4} + d_0 d_2 [(d_1 - 1)\kappa_1 + \frac{1}{2}] + \frac{\|\Phi\|_F^2}{4}, \quad (39)$$

$$\mathbb{E}[(G_1 G_1^\top)_{ac}] = d_0 d_2 \kappa_2 + O(d_2 \sqrt{d_0}) \quad (a \neq c). \quad (40)$$

As $d_0 \rightarrow \infty$, $\kappa_1 \rightarrow \frac{1}{4\pi^2} + \frac{1}{16}$ and $\kappa_2 \rightarrow \frac{1}{4\pi}$.

C.4. Preliminary lemmas

Lemma 16 Same as Lemma 1.

Lemma 17 (Gaussian integration by parts) Let $x \sim \mathcal{N}(0, \mathbf{I}_n)$ and $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be locally absolutely continuous along almost every line parallel to the j -th axis, with weak derivative $\partial_j f$, and $\mathbb{E}|x_j f(x)| < \infty$, $\mathbb{E}|\partial_j f(x)| < \infty$. Then

$$\mathbb{E}[x_j f(x)] = \mathbb{E}[\partial_j f(x)]. \quad (41)$$

In particular this holds whenever f is Lipschitz.

Proof. A standard identity [20, 25]. The standard Gaussian density φ satisfies $\partial_j \varphi(x) = -x_j \varphi(x)$, so by Fubini and one-dimensional integration by parts in x_j (the boundary term vanishes by integrability),

$$\mathbb{E}[x_j f(x)] = - \int f \partial_j \varphi = \int \partial_j f \varphi = \mathbb{E}[\partial_j f(x)].$$

The Lipschitz case follows from Rademacher's theorem. ■

Remark 2 $z \mapsto \sigma(z)$ is 1-Lipschitz, so $f(x) = \sigma(w^\top x)$ falls under Lemma 17 with $\partial_j f(x) = \mathbf{1}[w^\top x > 0] w_j$ a.e. By contrast $g(x) = \mathbf{1}[w^\top x > 0] h(x)$ with h not vanishing on $\{w^\top x = 0\}$ is discontinuous across that hyperplane and is not covered by Lemma 17; for such integrands we use a half-space argument (Lemma 21) or differentiate a separately established Lipschitz identity (Lemma 19).

Lemma 18 For $Z \sim \mathcal{N}(0, s^2)$ with $s > 0$,

$$\mathbb{E}[\sigma(Z)] = \frac{s}{\sqrt{2\pi}}, \quad \mathbb{E}[\sigma(Z)^2] = \frac{s^2}{2}, \quad \mathbb{E}[\sigma'(Z)] = \frac{1}{2}. \quad (42)$$

Proof. Direct computation: $\mathbb{E}[\sigma(Z)] = \int_0^\infty z \varphi_s(z) dz = s/\sqrt{2\pi}$; $\mathbb{E}[\sigma(Z)^2] = \int_0^\infty z^2 \varphi_s(z) dz = s^2/2$; $\mathbb{E}[\sigma'(Z)] = \Pr(Z > 0) = 1/2$. ■

Lemma 19 (Cho–Saul kernel) *Let $u, v \in \mathbb{R}^n \setminus \{0\}$, $x \sim \mathcal{N}(0, \mathbf{I}_n)$, $\theta := \angle(u, v)$. Then [5]*

$$\mathbb{E}[\sigma(u^\top x)\sigma(v^\top x)] = \frac{\|u\| \|v\|}{2\pi} (\sin \theta + (\pi - \theta) \cos \theta), \quad (43)$$

$$\mathbb{E}[\sigma'(u^\top x)\sigma'(v^\top x)] = \frac{\pi - \theta}{2\pi}. \quad (44)$$

Proof. Both quantities depend on $(u^\top x, v^\top x)$, bivariate Gaussian with variances $\|u\|^2, \|v\|^2$ and correlation $\cos \theta$. Equation (44) is $\Pr(u^\top x > 0, v^\top x > 0)$, the angular measure of the cone $\{Z_1 > 0, Z_2 > 0\}$ for a standard bivariate Gaussian with correlation $\cos \theta$, namely $(\pi - \theta)/(2\pi)$. For (43), by Price’s theorem [22],

$$\frac{\partial}{\partial \rho} \mathbb{E}[\sigma(u^\top x)\sigma(v^\top x)] = \|u\| \|v\| \mathbb{E}[\sigma'(u^\top x)\sigma'(v^\top x)] = \frac{\|u\| \|v\| (\pi - \theta)}{2\pi};$$

at $\rho = -1$ (i.e. $\theta = \pi$) the product vanishes a.s., and integrating ∂_ρ from -1 to $\cos \theta$ gives (43). ■

Lemma 20 (Wick/Isserlis) *Same as Lemma 3.*

Lemma 21 *Let $x \sim \mathcal{N}(0, \mathbf{I}_n)$ and $u \in \mathbb{R}^n \setminus \{0\}$. Then*

$$\mathbb{E}[xx^\top \mathbf{1}[u^\top x > 0]] = \frac{1}{2} \mathbf{I}_n. \quad (45)$$

Proof. Write $A := \mathbb{E}[xx^\top \mathbf{1}[u^\top x > 0]]$ and $x = (u^\top x / \|u\|)u + x_\perp$ with x_\perp the projection onto u^\perp . In an orthonormal basis containing $u/\|u\|$, the off-diagonal entries of A vanish by the reflection symmetry of x_\perp ; the entry along $u/\|u\|$ is $\mathbb{E}[(u^\top x / \|u\|)^2 \mathbf{1}[u^\top x > 0]] = \frac{1}{2}$ ($u^\top x$ symmetric), and each remaining diagonal entry is $\frac{1}{2} \mathbb{E}[x_j^2] = \frac{1}{2}$ by independence of x_\perp from $u^\top x$ and $\Pr(u^\top x > 0) = \frac{1}{2}$. Hence $A = \frac{1}{2} \mathbf{I}_n$. ■

C.5. Proof of the gradient self outer-product of W_2

Lemma 22 $G_2 = W_2 H - \Phi C^\top$, where $H := \mathbb{E}_x[h(x)h(x)^\top] \in \mathbb{R}^{d_1 \times d_1}$ and $C := \mathbb{E}_x[h(x)x^\top] \in \mathbb{R}^{d_1 \times d_0}$.

Proof. Differentiating (35) in W_2 ,

$$G_2 = \mathbb{E}_x[(W_2 h(x) - \Phi x)h(x)^\top] = W_2 \mathbb{E}_x[h(x)h(x)^\top] - \Phi \mathbb{E}_x[xh(x)^\top] = W_2 H - \Phi C^\top. \quad \blacksquare$$

Lemma 23 $C = \frac{1}{2} W_1$.

Proof. $C_{aj} = \mathbb{E}_x[\sigma(w_a^\top x)x_j]$, and $f(x) := \sigma(w_a^\top x)$ is Lipschitz with $\partial_j f(x) = \mathbf{1}[w_a^\top x > 0](w_a)_j$ a.e. By Lemma 17 and Lemma 18,

$$C_{aj} = \mathbb{E}_x[\partial_j f(x)] = (w_a)_j \mathbb{E}_x[\sigma'(w_a^\top x)] = \frac{1}{2}(w_a)_j.$$

■

Proof of (38). By Lemma 22,

$$G_2 G_2^\top = W_2 H^2 W_2^\top - W_2 H C \Phi^\top - \Phi C^\top H W_2^\top + \Phi C^\top C \Phi^\top.$$

Since $\mathbb{E}[W_2] = 0$, the cross terms vanish under $\mathbb{E}_{W_2}[\cdot \mid W_1]$, and by Lemma 16 (left identity, $M = H^2$),

$$\mathbb{E}[G_2 G_2^\top \mid W_1] = \text{tr}(H^2) \mathbf{I}_{d_2} + \Phi C^\top C \Phi^\top. \quad (46)$$

By Lemma 23, $C^\top C = \frac{1}{4} W_1^\top W_1$, and Lemma 16 (right identity, $N = \mathbf{I}_{d_1}$) gives $\mathbb{E}_{W_1}[W_1^\top W_1] = d_1 \mathbf{I}_{d_0}$, so $\mathbb{E}_{W_1}[C^\top C] = \frac{d_1}{4} \mathbf{I}_{d_0}$. Taking \mathbb{E}_{W_1} of (46),

$$\mathbb{E}[G_2 G_2^\top] = \mathbb{E}_{W_1}[\text{tr}(H^2)] \mathbf{I}_{d_2} + \frac{d_1}{4} \Phi \Phi^\top = \alpha \mathbf{I}_{d_2} + \frac{d_1}{4} \Phi \Phi^\top.$$

For the order of α , write $\text{tr}(H^2) = \sum_a H_{aa}^2 + \sum_{a \neq b} H_{ab}^2$. By Lemma 18, $H_{aa} = \frac{1}{2} \|w_a\|^2$ with $\|w_a\|^2 \sim \chi_{d_0}^2$, so

$$\sum_a \mathbb{E}[H_{aa}^2] = \frac{d_1 d_0 (d_0 + 2)}{4} = \Theta(d_0^2 d_1).$$

By (43), $H_{ab} = \|w_a\| \|w_b\| g(\theta_{ab})$ with $g(\theta) = (\sin \theta + (\pi - \theta) \cos \theta)/(2\pi)$ and $\theta_{ab} = \angle(w_a, w_b)$ independent of $\|w_a\|, \|w_b\|$, so $\mathbb{E}[H_{ab}^2] = d_0^2 \mathbb{E}[g(\theta)^2]$; since $\mathbb{E}[g(\theta)^2] \rightarrow g(\pi/2)^2 = 1/(4\pi^2) > 0$,

$$\sum_{a \neq b} \mathbb{E}[H_{ab}^2] = d_1 (d_1 - 1) d_0^2 \mathbb{E}[g(\theta)^2] = \Theta(d_0^2 d_1^2).$$

The off-diagonal term dominates, so $\alpha = \Theta(d_0^2 d_1^2)$. ■

C.6. Proof of the gradient self outer-product of W_1

By the chain rule, with residual $r(x) := W_2 \sigma(W_1 x) - \Phi x$ and $D(x) := \text{diag}(\sigma'(W_1 x))$,

$$G_1 = G_1^{(1)} - G_1^{(2)}, \quad G_1^{(1)} := \mathbb{E}_x[D(x) W_2^\top W_2 \sigma(W_1 x) x^\top], \quad G_1^{(2)} := \mathbb{E}_x[D(x) W_2^\top \Phi x x^\top]. \quad (47)$$

$G_1^{(1)}$ is quadratic in W_2 and $G_1^{(2)}$ is linear, so the cross terms in $G_1 G_1^\top$ are odd in W_2 and vanish under $\mathbb{E}_{W_2}[\cdot \mid W_1]$:

$$\mathbb{E}_{W_2}[G_1 G_1^\top \mid W_1] = \mathbb{E}_{W_2}[G_1^{(1)} (G_1^{(1)})^\top \mid W_1] + \mathbb{E}_{W_2}[G_1^{(2)} (G_1^{(2)})^\top \mid W_1]. \quad (48)$$

Lemma 24 $G_1^{(2)} = \frac{1}{2} W_2^\top \Phi$, and $\mathbb{E}_{W_2}[G_1^{(2)} (G_1^{(2)})^\top \mid W_1] = \frac{\|\Phi\|_F^2}{4} \mathbf{I}_{d_1}$.

Proof. Let $v_a := \Phi^\top(W_2)_{:,a} \in \mathbb{R}^{d_0}$, so $(W_2^\top \Phi)_{ab} = (v_a)_b$. Then, by Lemma 21,

$$(G_1^{(2)})_{ab} = \mathbb{E}_x[\mathbf{1}[w_a^\top x > 0] x_b (x^\top v_a)] = (\mathbb{E}_x[\mathbf{1}[w_a^\top x > 0] x x^\top] v_a)_b = \frac{1}{2} (v_a)_b,$$

so $G_1^{(2)} = \frac{1}{2} W_2^\top \Phi$ and $G_1^{(2)} (G_1^{(2)})^\top = \frac{1}{4} W_2^\top \Phi \Phi^\top W_2$; Lemma 16 (right identity, $N = \Phi \Phi^\top$) gives $\mathbb{E}_{W_2}[W_2^\top \Phi \Phi^\top W_2] = \|\Phi\|_F^2 \mathbf{I}_{d_1}$. ■

Lemma 25 Define $F_b(u, v) := \mathbb{E}_x[\sigma'(u^\top x)\sigma(v^\top x)x_b]$. Then $F_b(u, u) = \frac{1}{2}u_b$, and for $u \notin \text{span}(v)$,

$$F_b(u, v) = u_b \cdot \frac{\|v\| \sin \theta}{2\pi \|u\|} + v_b \cdot \frac{\pi - \theta}{2\pi}, \quad \theta := \angle(u, v). \quad (49)$$

Proof. For $u = v$, $\sigma'(u^\top x)\sigma(u^\top x) = \sigma(u^\top x)$ a.e., so by Lemma 17 (applied to the Lipschitz $x \mapsto \sigma(u^\top x)$) and Lemma 18, $F_b(u, u) = \mathbb{E}_x[\sigma(u^\top x)x_b] = \mathbb{E}_x[\sigma'(u^\top x)]u_b = \frac{1}{2}u_b$. For $u \neq v$, $f(x) = \sigma'(u^\top x)\sigma(v^\top x)$ is discontinuous across $\{u^\top x = 0\}$, so we differentiate the Cho–Saul kernel instead. With $K(u, v) := \mathbb{E}_x[\sigma(u^\top x)\sigma(v^\top x)]$, dominated convergence (since σ is Lipschitz with linear growth) gives

$$\frac{\partial K}{\partial u_b}(u, v) = \mathbb{E}_x[\sigma'(u^\top x)x_b\sigma(v^\top x)] = F_b(u, v).$$

By Lemma 19, $K(u, v) = \frac{\|u\|\|v\|}{2\pi}(\sin \theta + (\pi - \theta)\cos \theta)$, smooth for $u \notin \text{span}(v)$; differentiating with $\partial \|u\| / \partial u_b = u_b / \|u\|$, $\partial \cos \theta / \partial u_b = \frac{1}{\|u\|}(v_b / \|v\| - \cos \theta u_b / \|u\|)$, $\partial \theta / \partial u_b = -\frac{1}{\sin \theta} \partial \cos \theta / \partial u_b$, and simplifying yields (49). \blacksquare

Lemma 26 $\mathbb{E}_{W_2}[(W_2^\top W_2)_{ak}(W_2^\top W_2)_{cl}] = d_2^2 \delta_{ak}\delta_{cl} + d_2 \delta_{ac}\delta_{kl} + d_2 \delta_{al}\delta_{ck}$.

Proof. $(W_2^\top W_2)_{ak} = \sum_{i=1}^{d_2} (W_2)_{ia}(W_2)_{ik}$, and applying Lemma 20 to the four Gaussian factors $(W_2)_{ia}, (W_2)_{ik}, (W_2)_{jc}, (W_2)_{jl}$ (summed over i, j) the three pairings give $\delta_{ak}\delta_{cl}d_2^2 + \delta_{ac}\delta_{kl}d_2 + \delta_{al}\delta_{ck}d_2$. \blacksquare

Proof of (39)–(40). Permutation invariance follows from the exchangeability of the rows of W_1 and columns of W_2 . The $G_1^{(2)}$ contribution to (48) is $\frac{\|\Phi\|_F^2}{4} \mathbf{I}_{d_1}$ by Lemma 24. For $G_1^{(1)}$, write $(G_1^{(1)})_{ab} = \sum_{k=1}^{d_1} (W_2^\top W_2)_{ak} F_b(w_a, w_k)$, so

$$\mathbb{E}_{W_2}[(G_1^{(1)}(G_1^{(1)})^\top)_{ac} \mid W_1] = \sum_{b,k,l} \mathbb{E}_{W_2}[(W_2^\top W_2)_{ak}(W_2^\top W_2)_{cl}] F_b(w_a, w_k) F_b(w_c, w_l),$$

and substituting Lemma 26 gives $T_1^{(ac)} + T_2^{(ac)} + T_3^{(ac)}$ with

$$T_1^{(ac)} = d_2^2 \sum_b F_b(w_a, w_a) F_b(w_c, w_c), \quad (50)$$

$$T_2^{(ac)} = d_2 \delta_{ac} \sum_{k=1}^{d_1} \sum_b F_b(w_a, w_k)^2, \quad (51)$$

$$T_3^{(ac)} = d_2 \sum_b F_b(w_a, w_c) F_b(w_c, w_a). \quad (52)$$

Term T_1 . By Lemma 25, $F_b(w_a, w_a) = \frac{1}{2}(w_a)_b$, so $T_1^{(ac)} = \frac{d_2^2}{4} w_a^\top w_c$, giving $\mathbb{E}[T_1^{(aa)}] = \frac{d_0 d_2^2}{4}$ and $\mathbb{E}[T_1^{(ac)}] = 0$ for $a \neq c$.

Term T_2 . For $u \neq v$, writing $A := \frac{\|v\| \sin \theta}{2\pi \|u\|}$, $B := \frac{\pi - \theta}{2\pi}$ in (49) and using $u^\top v = \|u\| \|v\| \cos \theta$,

$$\sum_b F_b(u, v)^2 = A^2 \|u\|^2 + 2AB(u^\top v) + B^2 \|v\|^2 = \|v\|^2 \psi(\theta), \quad \psi(\theta) := \frac{\sin^2 \theta + (\pi - \theta)^2}{4\pi^2} + \frac{\sin \theta \cos \theta (\pi - \theta)}{2\pi^2}.$$

In $T_2^{(aa)}$ the term $k = a$ contributes $\sum_b F_b(w_a, w_a)^2 = \frac{1}{4} \|w_a\|^2$ (expectation $\frac{d_0}{4}$), and each $k \neq a$ contributes $\mathbb{E}[\|w_k\|^2] \mathbb{E}[\psi(\theta_{ak})] = d_0 \kappa_1$, so

$$\mathbb{E}[T_2^{(aa)}] = d_2 \left[\frac{d_0}{4} + (d_1 - 1)d_0 \kappa_1 \right] = d_0 d_2 \left[(d_1 - 1)\kappa_1 + \frac{1}{4} \right],$$

and $T_2^{(ac)} = 0$ for $a \neq c$.

Term T_3 . For $a = c$, $\sum_b F_b(w_a, w_a)^2 = \frac{1}{4} \|w_a\|^2$, so $\mathbb{E}[T_3^{(aa)}] = \frac{d_0 d_2}{4}$. For $a \neq c$, expanding $F_b(w_a, w_c) F_b(w_c, w_a)$ via (49) (with $\theta := \theta_{ac}$) and summing over b ,

$$\sum_b F_b(w_a, w_c) F_b(w_c, w_a) = \frac{\|w_a\| \|w_c\|}{4\pi^2} \left[\cos \theta (\sin^2 \theta + (\pi - \theta)^2) + 2 \sin \theta (\pi - \theta) \right].$$

Here $\|w_a\| \|w_c\|$ is independent of θ_{ac} with $\mathbb{E}[\|w_a\| \|w_c\|] = d_0 + O(1)$; as $d_0 \rightarrow \infty$, θ_{ac} concentrates at $\pi/2$ with $\mathbb{E}[\cos \theta_{ac}] = 0$ and $\text{Var}(\cos \theta_{ac}) = 1/d_0$, so $|\mathbb{E}[\cos \theta_{ac} q(\theta_{ac})]| = O(1/\sqrt{d_0})$ for bounded q , while $\mathbb{E}[\sin \theta_{ac} (\pi - \theta_{ac})] = \Theta(1)$. The $\cos \theta(\cdot)$ part contributes $O(d_2 \sqrt{d_0})$ and the $2 \sin \theta (\pi - \theta)$ part gives $d_0 d_2 \kappa_2 + O(d_2)$, so $\mathbb{E}[T_3^{(ac)}] = d_0 d_2 \kappa_2 + O(d_2 \sqrt{d_0})$.

Combining. For $a = c$, adding the three terms and the $G_1^{(2)}$ contribution and merging the two $\frac{d_0 d_2}{4}$ pieces,

$$\mathbb{E}[(G_1 G_1^\top)_{aa}] = \frac{d_0 d_2^2}{4} + d_0 d_2 \left[(d_1 - 1)\kappa_1 + \frac{1}{2} \right] + \frac{\|\Phi\|_F^2}{4};$$

for $a \neq c$ only $T_3^{(ac)}$ survives, $\mathbb{E}[(G_1 G_1^\top)_{ac}] = d_0 d_2 \kappa_2 + O(d_2 \sqrt{d_0})$. Finally $\psi(\pi/2) = \frac{1 + (\pi/2)^2}{4\pi^2} = \frac{1}{4\pi^2} + \frac{1}{16}$ and $\sin(\pi/2)(\pi - \pi/2)/(2\pi^2) = \frac{1}{4\pi}$ give the high-dimensional limits of κ_1, κ_2 . ■

C.7. Numerical verification

End-to-end simulations match the formulas above without analytical shortcuts: sample W_1, W_2 with i.i.d. $\mathcal{N}(0, 1)$ entries and a large batch X with $x_i \sim \mathcal{N}(0, \mathbf{I}_{d_0})$, compute G_1, G_2 by automatic differentiation on the empirical loss, average $G_i G_i^\top$ over many initializations, and estimate $\alpha, \kappa_1, \kappa_2$ by Monte Carlo (never via $C = \frac{1}{2} W_1, G_1^{(2)} = \frac{1}{2} W_2^\top \Phi$, or the Cho–Saul kernel). At $d_0 = 25, d_1 = 20, d_2 = 12, \Phi_{ij} \sim \mathcal{N}(0, 9)$ (so $\|\Phi\|_F^2 = 3211.10$), $n_{\text{init}} = 4000, n_x = 8000$:

Quantity	Simulation	Theory	Ratio
G_2 diagonal mean	11541.18	11549.25	0.999
G_2 off-diagonal projection coefficient \hat{c}	4.50	5.00	0.900
G_1 diagonal mean	2372.19	2364.78	1.003
G_1 off-diagonal mean	24.64	23.38	1.054

Figure 2 sweeps one of d_0, d_1, d_2 at a time; the empirical slopes track the predicted exponents (for instance G_1 off-diagonal vs d_1 has slope ≈ 0 , vs d_0 and d_2 slope ≈ 1 , while G_1 diagonal vs d_1 is sublinear at moderate d_1 because the $\frac{d_0 d_2^2}{4} + \frac{\|\Phi\|_F^2}{4}$ constants still dominate the $d_0 d_1 d_2 \kappa_1$ term).

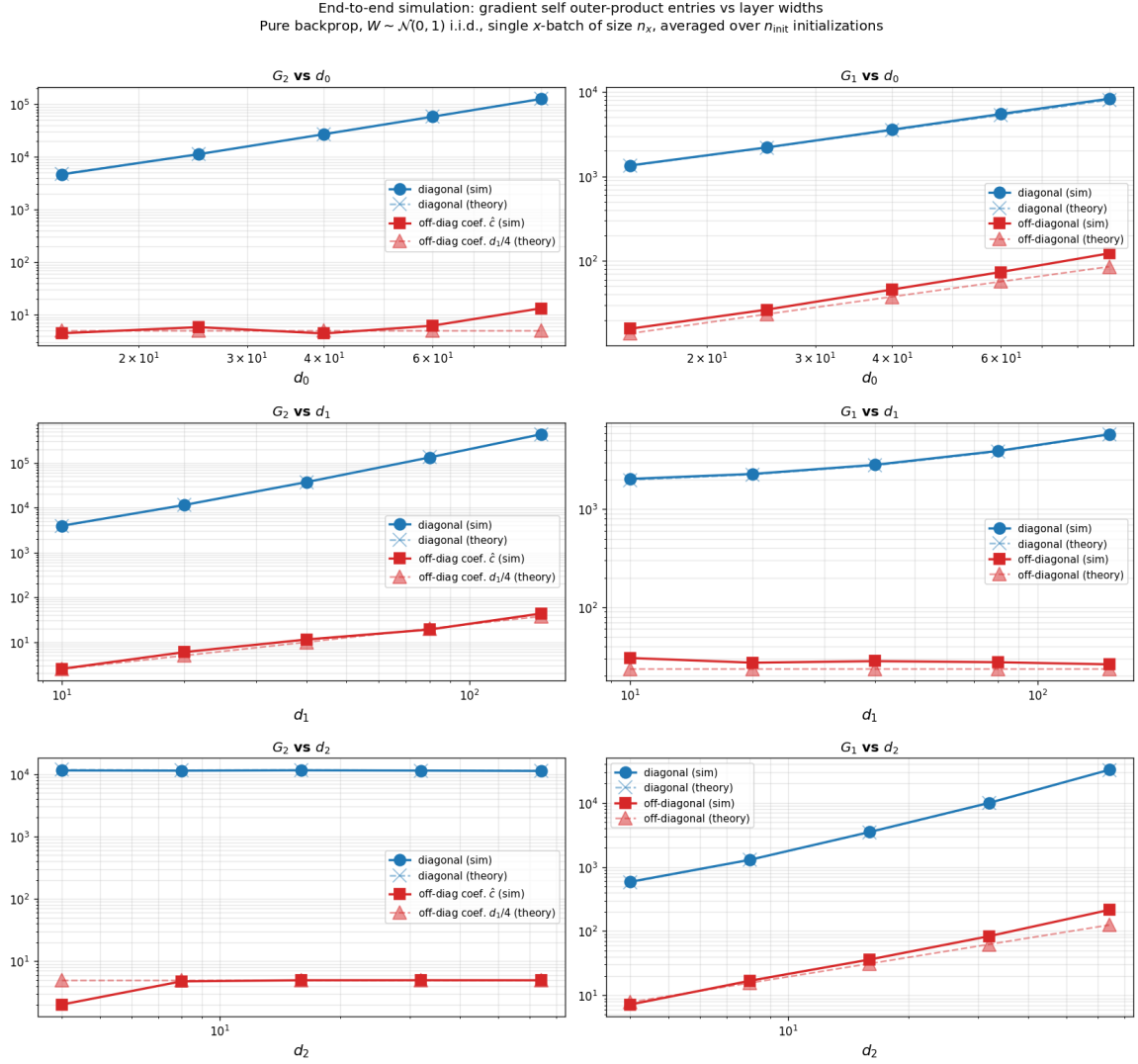


Figure 2: End-to-end simulation of the two-layer ReLU network. Rows: vary d_0 ($d_1=20, d_2=12$); vary d_1 ($d_0=25, d_2=12$); vary d_2 ($d_0=25, d_1=20$). Left: G_2 ; right: G_1 . Diagonal entries (circles, simulation) match theory (crosses); off-diagonal entries (squares, simulation) match theory (triangles).

C.8. Discussion

At initialization the gradient of W_2 is mostly target-independent: the $\Theta(d_0^2 d_1^2)$ leading order of $\mathbb{E}[(G_2 G_2^\top)_{aa}]$ comes entirely from $\text{tr}(H^2)$, which depends only on W_1 , while the target term $\frac{d_1}{4} (\Phi \Phi^\top)_{aa}$ is dominated unless $\|\Phi\|_F^2$ is exceptionally large. The nonzero $\Theta(d_0 d_2)$ off-diagonal of $\mathbb{E}[G_1 G_1^\top]$ comes from T_3 and reflects ReLU’s lack of odd symmetry: in a deep linear network the symmetry $W_1 \mapsto DW_1, W_2 \mapsto W_2 D$ for $D = \text{diag}(\pm 1)$ preserves the loss and forces off-diagonal

expectations to zero (cf. Section B), whereas $\sigma(-z) \neq -\sigma(z)$ breaks it. The dominance is therefore asymptotic in d_1 rather than exact.

Appendix D. Orthogonalization Is Asymptotically a Row Normalization

What does diagonal dominance say about the orthogonalization step itself? Suppose $\mathbf{V}_t \mathbf{V}_t^\top$ is diagonal. Then $\mathbf{V}_t \mathbf{V}_t^\top = \text{diag}(\mathbf{V}_t \mathbf{V}_t^\top)$, so $(\mathbf{V}_t \mathbf{V}_t^\top)^{-1/2} = \text{diag}(\mathbf{V}_t \mathbf{V}_t^\top)^{-1/2}$, and since $\text{diag}(\mathbf{V}_t \mathbf{V}_t^\top)_{ii} = \|(\mathbf{V}_t)_{i,:}\|_2^2$,

$$[(\mathbf{V}_t \mathbf{V}_t^\top)^{-1/2} \mathbf{V}_t]_{i,:} = \frac{(\mathbf{V}_t)_{i,:}}{\|(\mathbf{V}_t)_{i,:}\|_2}. \quad (53)$$

That is, when $\mathbf{V}_t \mathbf{V}_t^\top$ is diagonal the orthogonalization in (1) is just each row of the momentum divided by its own ℓ_2 norm, which is exactly the row normalization used by several recent stateless optimizers [7, 10, 18, 21, 23, 29]; equivalently, its implicit preconditioner is $\text{diag}(\mathbf{P}_t)^{1/2} \otimes \mathbf{I}_n$, which equals H_{Muon} when \mathbf{P}_t is diagonal. Putting (53) together with Theorems 1 to 3: at a Gaussian initialization, as the width grows, Muon’s $\Theta(mn \min(m, n))$ orthogonalization collapses to a $\Theta(mn)$ row-wise ℓ_2 normalization with the same preconditioner shape. Deng et al. [7] check that this equivalence holds throughout the training of GPT-2 and LLaMA, and that the dominance ratio grows with model size, in line with the width dependence in Theorem 3. Carrying the analysis past initialization, and to deep nonlinear networks, is left for future work.