
Low-Rank Neural Networks and Finite-Width NTK at the Edge of Convexity

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Low-rank neural networks are usually sold as smaller, compressed and pruned
2 models. This paper asks a sharper question: when does low rank still preserve
3 the optimization geometry that makes wide networks trainable? We answer this
4 through the finite-width neural tangent kernel. Low-rank NTK training still carries
5 an interpretable bottleneck-feature geometry, because rank controls which feature
6 maps enter the tangent kernel. Our main result is that low-rank training is governed
7 by a genuine three-way compromise between width, depth, and rank. To the best
8 of our knowledge, this is the first NTK analysis showing that an optimization
9 certificate can impose a depth-rank tradeoff of this form. In particular, for low-
10 rank bottleneck dynamics, preserving the NTK spectral margin has a proven
11 conservative cubic-depth sufficient rule, $r \gtrsim L^3$, up to dataset-size and logarithmic
12 factors. From the finite-network experiments, we conjecture an effective $L^{3/2}$
13 law for scalar-output contracted cumulants. We also give a full-rank-compatible
14 parametrization, separating true low-rank effects from ordinary finite-width effects.
15 True finite-network NTK experiments confirm the exact full-rank matching, the
16 predicted operator scaling, the rank-depth tradeoff, and the finite-network cumulant
17 growth. The result is a criterion for when low-rank networks remain trainable for
18 structural reasons, not merely parameter-efficient ones.

19 1 Introduction

20 Low-rank neural networks are usually motivated by efficiency: fewer parameters, lower memory
21 traffic, and cheaper adaptation. This motivation is incomplete. Rank also changes the geometry of
22 training. Low-rank bottlenecks do not merely shrink a dense model: they restrict the network to a
23 small set of bottleneck feature maps, so feature selection already enters through the architecture. Thus
24 saying that an NTK analysis has “no features” is misleading for low-rank networks; the question is
25 whether the fixed or linearized low-rank feature geometry remains stable enough to explain training.
26 In the NTK regime, gradient descent is close to kernel regression when the empirical NTK Gram
27 matrix is positive and stable. Thus a low-rank network is useful not only when it is small, but when
28 its finite NTK stays close enough to a positive limiting kernel. We call this the edge of convexity: the
29 regime where the network is finite and structured, but the NTK Gram matrix still provides a convex
30 local training geometry.

31 Empirically, low-rank bottlenecks can give a much better parameter-accuracy tradeoff than dense
32 baselines. This raises the next theoretical question at very small rank: if only r bottleneck feature
33 maps are available, which features are recovered, and is the recovery explained by a kernel/NTK
34 mechanism or by genuine feature-learning dynamics such as mean-field or μ P scaling? The present
35 paper answers the controlled NTK side of that question. It shows that even in the nominally lazy

Table 1: MNIST parameter-efficiency snapshot. Low-rank bottlenecks recover most of the dense baseline accuracy with far fewer trainable parameters.

Model	Trainable parameters	Test accuracy (%)
Dense MLP	669,706	98.39
Low-rank, $r = 5$	7,695	93.72
Low-rank, $r = 10$	10,260	96.24
Low-rank, $r = 15$	12,825	96.97
Low-rank, $r = 25$	17,955	97.00
Low-rank, $r = 50$	30,780	97.01
Low-rank trainable factors, $r = 32$	440,362	98.30

36 regime, low-rank bottlenecks introduce nontrivial rank-channel geometry, memory suppression, and
 37 spectral constraints that any feature-learning theory must build on.

38 This question is difficult because three limits interact. Width controls ordinary finite-width fluctua-
 39 tions, depth amplifies NTK tensor corrections, and rank changes the channel contractions themselves.
 40 Treating low rank as merely replacing $1/n$ by $1/r$ misses the full-rank endpoint: if a rank- n fac-
 41 torization is just an orthogonal change of variables, the low-rank correction must vanish exactly.
 42 Conversely, bottleneck random-feature low-rank networks do not reduce to dense MLPs and can
 43 suppress old NTK paths geometrically. A useful theory must separate these two cases.

44 We study random-feature low-rank (RF-LR) scalar-output ReLU networks with layer weights of the
 45 form

$$W^{(\ell)} = \frac{\sigma_w}{\sqrt{n_{\ell-1}}} \sqrt{\frac{n_{\ell}}{r_{\ell}}} U^{(\ell)} B^{(\ell)}, \quad (U^{(\ell)})^{\top} U^{(\ell)} = I_{r_{\ell}}, \quad (1)$$

46 where $U^{(\ell)}$ is frozen and sampled uniformly from the Stiefel manifold $\text{St}(n_{\ell}, r_{\ell}) = \{U \in \mathbb{R}^{n_{\ell} \times r_{\ell}} : U^{\top} U = I_{r_{\ell}}\}$, while $B^{(\ell)}$ is trainable. Thus the Stiefel object is the left factor $U^{(\ell)}$ itself: its columns
 47 form an orthonormal r_{ℓ} -frame in $\mathbb{R}^{n_{\ell}}$, chosen uniformly among all such frames. Equivalently, $U^{(\ell)}$
 48 selects a random r_{ℓ} -dimensional subspace with no preferred direction. When $r_{\ell} = n_{\ell}$, the map
 49 $B^{(\ell)} \mapsto U^{(\ell)} B^{(\ell)}$ is an orthogonal reparametrization of a dense MLP. Thus the full-rank endpoint is
 50 exact, not asymptotic; Section 4 identifies the corresponding rank-defect scale.
 51

52 2 Related Work

53 The NTK formalizes lazy infinite-width training as kernel regression [1, 2]. Lee et al. clarified that
 54 wide networks of any depth evolve approximately as linear models around initialization, making the
 55 empirical NTK spectrum the relevant stability object [2]. Many global-convergence results rely on
 56 keeping this empirical NTK close to a positive limiting kernel [3, 4, 5]. Spectral refinements show
 57 that strong guarantees can hold even when only part of a deep network is very wide [6, 7]; recent EOC
 58 analyses further study concentration and spectra of MLP NTKs at the edge of chaos [8, 9]. These
 59 works provide the dense full-rank baseline; they do not identify the finite-rank correction created by
 60 factorized layers or the cancellation that occurs at the full-rank endpoint.

61 A complementary line studies how finite networks depart from their infinite-width limits. Yang’s
 62 Tensor Programs give a general language for infinite-width limits, parametrization choices, and
 63 feature-learning scalings across architectures [10, 11, 12, 13], while the neural tangent hierarchy
 64 describes finite-width departures from frozen-kernel dynamics [18]. Feynman-diagram methods
 65 organize these corrections by connected cumulants and tensor contractions [19, 20]. Our analysis
 66 follows this perturbative viewpoint, but changes the stochastic source: dense channel contractions are
 67 replaced by Stiefel-projector and rank-channel cumulants.

68 NTK theory is also not meant to replace feature-learning theories. Mean-field limits describe
 69 distributional feature evolution beyond the fixed-kernel regime [15, 16, 17], and μP , depth- μP , tensor-
 70 program, and DMFT analyses identify scalings where features continue to move at infinite width
 71 [12, 13, 14]. Our results are complementary: they quantify the finite-rank and finite-width stability of
 72 the lazy/proxy geometry before one studies full feature-learning dynamics.

73 Low-rank adaptation and low-rank training have also been analyzed in NTK-like settings, including
 74 guarantees that sufficiently large low-rank adapters avoid bad local minima in certain regimes [22].
 75 Empirical NTK tools make direct finite-network checks feasible [21], and recent quantitative GP
 76 approximation during training gives another complementary control problem [23]. In contrast,
 77 our object is the operator deviation of a deep low-rank empirical NTK and the resulting Weyl
 78 spectral-preservation criterion.

79 **Contributions.** First, full-rank factorization is only a coordinate change: when $r = n$, the empirical
 80 NTK matches the dense MLP pathwise. Away from full rank, the correction is a projector defect, not
 81 the naive replacement $n \mapsto r$.

82 The second message is that depth turns small rank defects into a stability constraint. Low-rank
 83 randomness enters locally at each layer, but NTK legs are transported through the whole network.
 84 One transported leg creates a larger accumulated correction than a pure NNGP fluctuation, and two
 85 transported NTK legs create the worst depth growth. This is the mechanism behind the conservative
 86 cubic rank-depth rule for bottleneck dynamics: deep low-rank networks are hard to keep in the same
 87 convex NTK regime unless rank grows with depth.

88 The third message is practical. The theory gives an operator-level stability certificate under an
 89 explicit well-spread data assumption, and it remains closed under the NTK descendants that control
 90 short-time training drift. The experiments then check the pieces separately: exact full-rank matching,
 91 disappearance of the rank defect at $r = n$, the predicted depth normalization, and direct finite-network
 92 cumulants. The observed slopes are effective exponents of mixed observable quantities, so they also
 93 point to a sharper future theory beyond the worst-case cubic envelope.

94 3 Models and Full-Rank Matching

95 This section defines the factorized network, proves exact matching at full rank, and explains why the
 96 left factor is frozen.

97 For a fixed dataset x_1, \dots, x_m , the scalar-output empirical NTK is

$$\widehat{\Theta}_{ab} = \nabla_{\theta} f(x_a) \cdot \nabla_{\theta} f(x_b). \quad (2)$$

98 In (1), the trainable parameters are the right factors $B^{(\ell)}$ and the readout. We keep $U^{(\ell)}$ frozen so the
 99 tangent space is a fixed subspace and all rank-dependent randomness enters through $G = (n/r)UU^{\top}$.
 100 Training U would move the subspace and add U - B dNTK interaction terms.

101 **Theorem 3.1** (Pathwise full-rank matching). *If $r_{\ell} = n_{\ell}$ for every hidden layer, then for every*
 102 *realization, dataset, and depth, the empirical NTK in B -coordinates equals the dense MLP NTK*
 103 *in coordinates $\widehat{W}^{(\ell)} = U^{(\ell)}B^{(\ell)}$. The same equality holds for the NNGP, dNTK, ddNTK, and their*
 104 *tensor cumulants.*

105 *Proof.* When $r = n$, U is orthogonal, so $\Delta B \mapsto U\Delta B$ preserves Frobenius inner products. Hence
 106 $J_B J_B^{\top} = J_W J_W^{\top}$, and higher derivatives and cumulants are invariant under the same coordinate
 107 change. \square

108 **Proposition 3.2** (Why the left factor is frozen). *For one layer $W = \alpha UB$, with $U^{\top}U = I_r$ and*
 109 *$H = \partial\mathcal{L}/\partial W$, training only B gives*

$$\nabla_B \mathcal{L} = \alpha U^{\top} H, \quad \dot{W}_B = -\eta \alpha^2 U U^{\top} H. \quad (3)$$

110 *If U is also trained on the Stiefel manifold, the normal projected component contributes*

$$\dot{W}_U = -\eta \alpha^2 (I - U U^{\top}) H B^{\top} B \quad \text{up to rotations inside } \text{span}(U) \text{ that can be absorbed into } B. \quad (4)$$

111 *Thus moving U adds U - B interaction blocks proportional to $B^{\top} B$, absent from the frozen-left model*
 112 *and not controlled solely by cumulants of $G = (n/r)UU^{\top}$.*

113 *Proof idea.* Project the Euclidean gradient onto the B -directions and the Stiefel tangent space, then
 114 map both variations back to $W = \alpha UB$; this gives (3) and (4). The full differential calculation is in
 115 Appendix D. \square

116 Freezing U therefore isolates the rank-width effect and keeps the diagrammatics closed under the
 117 Stiefel projector vertices of Lemma 4.1.

118 The RF-LR bottleneck has a different recursion,

$$\Theta^{(\ell)} = 1 + \frac{1}{r} \Theta^{(\ell-1)} \dot{\Sigma}^{(\ell)} + \frac{1}{r} \Sigma^{(\ell)}. \quad (5)$$

119 Here we suppress the sample pair (x, x') . If $(g_x^{(\ell)}, g_{x'}^{(\ell)})$ is the Gaussian preactivation pair induced by
 120 the previous-layer covariance, then

$$\Sigma^{(\ell)}(x, x') := \mathbb{E} \left[\sigma(g_x^{(\ell)}) \sigma(g_{x'}^{(\ell)}) \right], \quad \dot{\Sigma}^{(\ell)}(x, x') := \mathbb{E} \left[\sigma'(g_x^{(\ell)}) \sigma'(g_{x'}^{(\ell)}) \right]. \quad (6)$$

121 Thus $\Sigma^{(\ell)}$ is the NNGP covariance term injected at layer ℓ , while $\dot{\Sigma}^{(\ell)}$ is the derivative covariance
 122 multiplying an old NTK line. At the ReLU fixed point $\dot{\Sigma}^{(\ell)} \rightarrow 1/2$, so an old NTK contribution is
 123 multiplied by approximately $(2r)^{-(L-k)}$. This is a different model, not a perturbation of the dense
 124 MLP.

125 4 Stiefel Projector Vertices

126 This section explains the random projector $G = (n/r)UU^\top$, which is the only way the frozen
 127 random left factor enters the NTK corrections. Let $U \in \mathbb{R}^{n \times r}$ be sampled uniformly from the Stiefel
 128 manifold and

$$G = \frac{n}{r} UU^\top. \quad (7)$$

129 We study this Stiefel projector because the trainable gradients with respect to B are always projected
 130 through the random subspace selected by U ; therefore cumulants of G are the low-rank replacement
 131 for dense neuron-index contractions. Orthogonal invariance gives the full covariance tensor.

132 **Lemma 4.1** (Stiefel projector covariance). *For all indices,*

$$\begin{aligned} \text{Cov}(G_{ij}, G_{kl}) &= \gamma_{n,r} \left(\delta_{ik} \delta_{jl} + \delta_{il} \delta_{jk} - \frac{2}{n} \delta_{ij} \delta_{kl} \right), \\ \gamma_{n,r} &= \frac{n(n-r)}{r(n-1)(n+2)} = \frac{1}{r} - \frac{1}{n} + O(n^{-2}). \end{aligned} \quad (8)$$

133 *In particular $\gamma_{n,n} = 0$. The mixed perturbation parameter is therefore $1/n + \gamma_{n,r}$: the $1/n$ term
 134 is the ordinary dense finite-width correction, while $\gamma_{n,r}$ is the rank defect. The rank-defect part
 135 vanishes at $r = n$, exactly as required by pathwise full-rank matching.*

136 *Proof.* The covariance of a Stiefel projector has the invariant form $a(\delta_{ik} \delta_{jl} + \delta_{il} \delta_{jk}) + b \delta_{ij} \delta_{kl}$. The
 137 identity $\text{tr } G = n$ gives $2a + nb = 0$. The standard second moment of a Stiefel projector gives
 138 $a = \gamma_{n,r}$, hence (8). \square

139 Higher cumulants are given by orthogonal Weingarten contractions. Connected s -point cumulants
 140 scale as $O_s(r^{1-s})$ and vanish in the full-rank endpoint. Thus the low-rank Feynman rule is simple:
 141 dense Gaussian-ReLU vertices remain, and each low-rank correction inserts connected cumulants of
 142 G .

143 The rank-channel power counting is the usual connected-cumulant counting for iid empirical averages:
 144 a connected s -point rank vertex scales as r^{1-s} . Lemma E.1 in Appendix E gives the precise statement
 145 and defines the channel atom $Z_a^{(\ell)}$. This is why finite-rank bias terms are typically $O(r^{-1})$, while the
 146 full-rank endpoint must be treated through the projector defect $G - I$, not by the informal replacement
 147 $n \mapsto r$.

148 4.1 Triangular tensor recursion

149 The full pairing-basis dictionary for V, D, F, A, B is in Appendix G.7. In the main text we use
 150 the compressed block notation $\mathcal{V} = (V_{1234}, V_{1324}, V_{1423})$, $\mathcal{C} = (D_{1234}, F_{1324}, F_{1423})$, and $\mathcal{A} =$
 151 $(A_{1234}, B_{1324}, B_{1423})$.

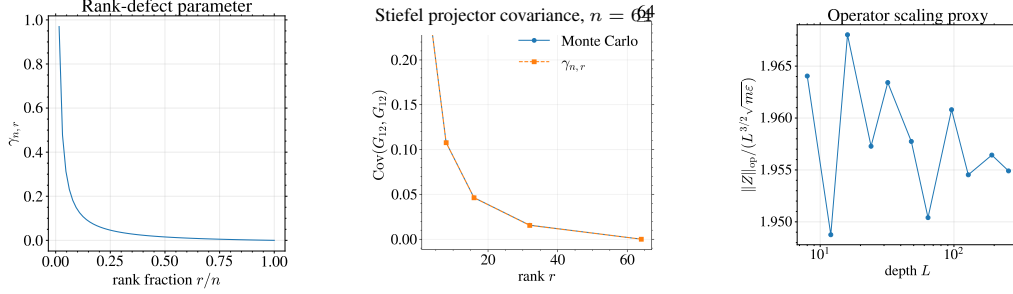


Figure 1: Rank-projector and operator checks. Left: the rank-defect scale $\gamma_{n,r}$ vanishes at full rank. Middle: empirical Stiefel-projector covariance matches the closed form. Right: the normalized operator proxy $\|Z\|_{\text{op}}/(L^{3/2}\sqrt{m\epsilon})$ is approximately stable across depth.

152 The layer recursion is triangular at first perturbative order:

$$\mathcal{V}_{\ell+1} = \mathcal{L}_V^{(\ell)} \mathcal{V}_\ell + \mathcal{S}_V^{(\ell)}, \quad (9)$$

$$\mathcal{C}_{\ell+1} = \mathcal{L}_C^{(\ell)} \mathcal{C}_\ell + \mathcal{B}_C^{(\ell)} \mathcal{V}_\ell + \mathcal{S}_C^{(\ell)}, \quad (10)$$

$$\mathcal{A}_{\ell+1} = \mathcal{L}_A^{(\ell)} \mathcal{A}_\ell + \mathcal{B}_A^{(\ell)} \mathcal{C}_\ell + \mathcal{Q}_A^{(\ell)} [\mathcal{V}_\ell, \mathcal{V}_\ell] + \mathcal{S}_A^{(\ell)}. \quad (11)$$

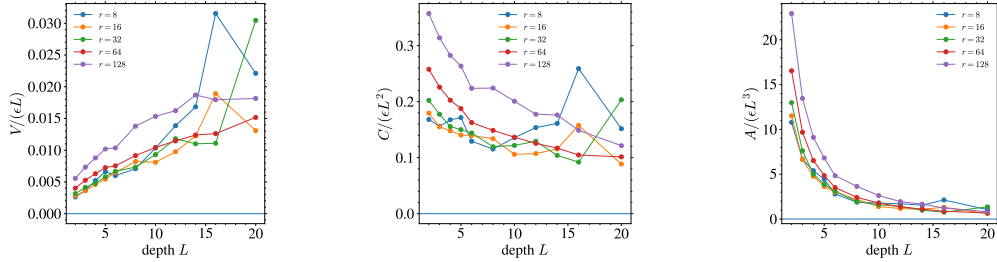


Figure 2: True finite-network cumulants from autodiff NTKs. These scalar-output cumulants are not isolated Feynman coefficients; they test the predicted polynomial control and rank-width scaling on actual finite networks.

153 The source bounds used in Proposition 5.2 are

$$\|\mathcal{S}_V^{(\ell)}\| \leq C, \quad \|\mathcal{S}_C^{(\ell)}\| \leq C(1+\ell+\|\mathcal{V}_\ell\|), \quad \|\mathcal{S}_A^{(\ell)}\| \leq C(1+\ell^2+\ell\|\mathcal{V}_\ell\|+\|\mathcal{C}_\ell\|+\|\mathcal{V}_\ell\|^2). \quad (12)$$

154 These formulas are the visible mechanism behind the hierarchy: V feeds D, F , and D, F feed A, B .

155 5 Tensor Scaling at ReLU Edge of Chaos

156 This section turns the projector cumulants into depth scalings for the tensor blocks V, D, F, A, B .
 157 For normalized ReLU, the edge-of-chaos correlation asymptotics used in the RF-LR analysis give
 158 $\theta_\ell = 3\pi/\ell + O(\log \ell/\ell^2)$, and

$$\dot{C}(\rho_\ell) = 1 - \frac{3}{\ell} + O\left(\frac{\log \ell}{\ell^2}\right). \quad (13)$$

159 Transporting an external NTK leg from layer k to L gives the factor $(k/L)^\lambda$ for an exponent
 160 $\lambda \in \{0, 3\}$, depending on whether the Gram entry is diagonal or off-diagonal.

161 **Lemma 5.1** (Depth-transport accumulation). *Let $T_\ell \in \mathbb{R}^d$, $\ell \geq 2$, satisfy $T_2 = 0$ and*

$$T_{\ell+1} = \left(I - \frac{M}{\ell} + E_\ell\right) T_\ell + \ell^\beta (s + e_\ell), \quad (14)$$

162 where $\|E_\ell\| \leq C \log^p(\ell)/\ell^2$, $\|e_\ell\| \leq C \log^q(\ell)/\ell$, every eigenvalue of M has real part greater than
 163 $-(\beta + 1)$, and $M + (\beta + 1)I$ is invertible. Then

$$T_L = L^{\beta+1} (M + (\beta + 1)I)^{-1} s + O(L^\beta \log^{p+q+2} L). \quad (15)$$

164 *Proof idea.* Iterating (14) gives a variation-of-constants formula; Appendix A gives the full discrete
 165 proof with the product estimates. The product of propagators from k to L is $(k/L)^M$ up to a
 166 summable perturbative error. Thus, at leading order,

$$T_L = L^{-M} \sum_{k \leq L} k^{M+\beta I} s + O(L^\beta \log^{p+q+2} L).$$

167 Euler-Maclaurin yields

$$L^{-M} \sum_{k \leq L} k^{M+\beta I} = L^{\beta+1} \int_0^1 u^{M+\beta I} du + O(L^\beta),$$

168 and the matrix integral is $(M + (\beta + 1)I)^{-1}$. \square

169 Let $\bar{\varepsilon} = \max_{\ell \leq L} (1/n_\ell + \gamma_{n_\ell, r_\ell})$. The finite-width tensors obey

$$V_L = O(\bar{\varepsilon}L), \quad C_L = O(\bar{\varepsilon}L^2), \quad A_L = O(\bar{\varepsilon}L^3), \quad (16)$$

170 where C denotes the mixed NNGP-NTK block (D, F) , and A denotes the NTK-NTK block (A, B) .

171 **Proposition 5.2** (Conditional tensor hierarchy). *Assume the low-rank cumulant recursion closes on*
 172 *the same tensor sectors as the dense finite-width Feynman hierarchy and that the ReLU edge-of-chaos*
 173 *transport is uniform away from the diagonal singularity. Then the isolated tensor coordinates satisfy*

$$V_L = O(\bar{\varepsilon}L), \quad D_L, F_L = O(\bar{\varepsilon}L^2), \quad A_L, B_L = O(\bar{\varepsilon}L^3). \quad (17)$$

174 *Consequently, the perturbative mean-NTK correction is controlled in the regime $L^3 \bar{\varepsilon} \ll 1$.*

175 *Proof idea.* Appendix A gives the complete block-recursion proof. The fresh NNGP four-point
 176 source contributes $O(\bar{\varepsilon})$ per layer and has no transported NTK leg in the dominant radial mode, so
 177 summation gives $V_L = O(\bar{\varepsilon}L)$. The mixed blocks D, F have one transported NTK leg and therefore
 178 one depth-transport accumulation with $\beta = 1$, yielding $O(\bar{\varepsilon}L^2)$ by Lemma 5.1. The NTK-NTK
 179 blocks A, B have two transported NTK legs, equivalently source exponent $\beta = 2$, giving $O(\bar{\varepsilon}L^3)$.
 180 These are isolated tensor coordinates; scalar-output cumulants can mix these sectors. \square

181 **Endpoint constants.** At the ReLU endpoint, set $X_0 = 2(g_+)^2$ and $Y_0 = 2\mathbf{1}_{\{g>0\}}$. Then

$$\text{Var}(X_0) = 5, \quad \text{Cov}(X_0, Y_0) = 1, \quad \text{Var}(Y_0) = 1. \quad (18)$$

182 The table should be read together with the plots. On the standard long run, the visible slopes
 183 are indeed close across ranks, with effective exponents around 1.5–2.1 for the same-off-diagonal
 184 observables. This supports the interpretation that the proven L^3 tensor bound is a safe upper envelope,
 185 not a claim that every scalar-output cumulant must scale exactly as L^3 . A natural conjecture is that
 186 some contracted observables obey a sharper mode-dependent law between $L^{3/2}$ and L^2 , while the
 187 cubic tensor hierarchy remains the conservative certificate. These are effective slopes over the plotted
 188 depth range, not asymptotic exponents.

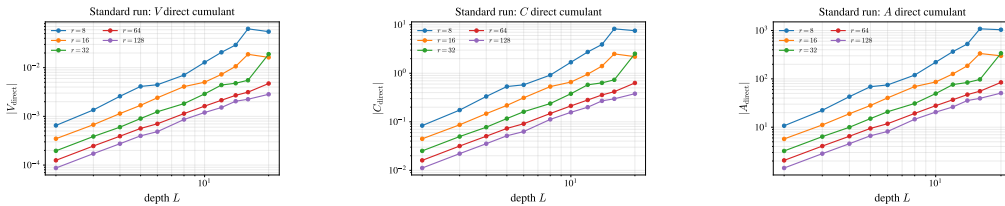


Figure 3: Direct finite-network cumulant magnitudes from the standard long run. From left to right: V_{direct} , C_{direct} , and A_{direct} . These scalar-output observables mix several diagrammatic sectors, so their slopes are effective diagnostics rather than isolated tensor exponents.

Table 2: Effective log-log depth slopes from the same standard long run as Figure 3, using $L \in \{2, 3, 4, 5, 6, 8, 10, 12, 14, 16, 20\}$. These slopes are diagnostics of scalar-output contracted observables, not isolated diagrammatic tensor exponents.

Rank r	V_{direct}	C_{direct}	A_{direct}
8	2.01	2.03	2.06
16	1.76	1.79	1.81
32	1.79	1.81	1.83
64	1.56	1.58	1.59
128	1.55	1.57	1.58

189 For a transported off-diagonal leg, $\lambda = 3$. Indeed, the leg is multiplied from layer k to layer L by

$$\prod_{j=k+1}^L \dot{C}(\rho_j) = \prod_{j=k+1}^L \left(1 - \frac{3}{j} + O\left(\frac{\log j}{j^2}\right) \right) = \left(\frac{k}{L}\right)^3 (1 + o(1)),$$

190 because the logarithm of the product is $-3 \sum_{j=k+1}^L j^{-1} + O(1/k)$. This is the source of the exponent
 191 3. Appendix K gives the ReLU moment and Riemann-sum calculation behind the mixed and
 192 NTK-NTK endpoint integrals:

$$B_\lambda = \frac{1}{\lambda+1} \left(5 - \frac{4}{\lambda+2} \right), \quad A_{\lambda,\mu} = \frac{1}{(\lambda+1)(\mu+1)} \left(5 - \frac{4}{\lambda+2} - \frac{4}{\mu+2} + \frac{4}{\lambda+\mu+3} \right). \quad (19)$$

193 Thus for off-diagonal modes,

$$V : 5, \quad C : B_3 = \frac{21}{20}, \quad A : A_{3,3} = \frac{173}{720}. \quad (20)$$

194 The diagonal/off-diagonal endpoint table is

Tensor coordinate	leading coefficient
$V(e, f)$	5
$D(e, f), F(e, f)$, transported leg diagonal	$B_0 = 3$
$D(e, f), F(e, f)$, transported leg off-diagonal	$B_3 = 21/20$
$A(e, f), B(e, f)$, both legs diagonal	$A_{0,0} = 7/3$
$A(e, f), B(e, f)$, one diagonal and one off-diagonal leg	$A_{0,3} = 43/60$
$A(e, f), B(e, f)$, both legs off-diagonal	$A_{3,3} = 173/720$

196 Equivalently, for four distinct generic inputs and the three pairing vectors in (51)–(53),

$$\mathcal{V}_L = 5 \bar{\varepsilon} L \mathbf{1}_3 + O(\bar{\varepsilon} \log^2 L), \quad (21)$$

$$\mathcal{C}_L = \frac{21}{20} \bar{\varepsilon} L^2 \mathbf{1}_3 + O(\bar{\varepsilon} L \log^2 L), \quad (22)$$

$$\mathcal{A}_L = \frac{173}{720} \bar{\varepsilon} L^3 \mathbf{1}_3 + O(\bar{\varepsilon} L^2 \log^2 L), \quad (23)$$

197 up to the finite pairing-basis projection factors already absorbed in the $O(\bar{\varepsilon})$ normalization. This is
 198 the calculation from which the main L, L^2, L^3 hierarchy is extracted.

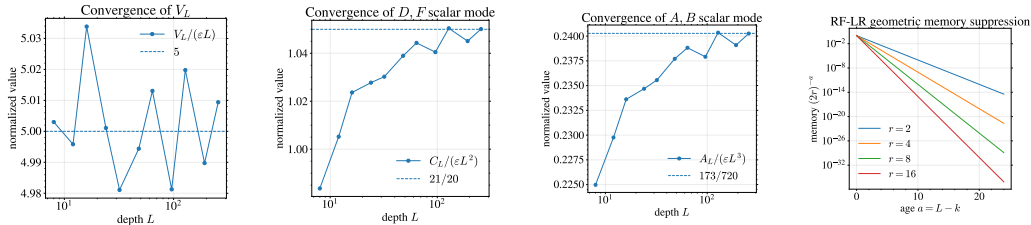


Figure 4: Tensor constants and RF-LR memory. The first three panels show convergence to the isolated constants 5, 21/20, and 173/720. The last panel shows geometric suppression of old NTK memory in the RF-LR bottleneck.

199 These constants are tensor coefficients after pairing-basis extraction, not arbitrary scalar-output cumu-
 200 lants. Direct finite-network cumulants mix Wick/Wishart terms, radial modes, readout contractions,
 201 and several Feynman sectors; Appendix H reports the corresponding log-log diagnostics.

202 A concrete three-layer delta-method calculation is given in Appendix B. It shows explicitly that
 203 right-factor-only low rank produces an r^{-1} mean correction and $r^{-1/2}$ single-seed fluctuations, while
 204 the diagonal ReLU correction cancels by scale invariance.

205 6 Operator Concentration and Spectral Preservation

206 The previous tensor estimates control individual NTK entries. For training stability, the relevant
 207 question is stronger: whether the whole empirical NTK matrix stays close enough to its limiting
 208 kernel that its smallest useful eigenvalue remains positive. Let $\Theta_{\infty,L}$ be the deterministic limiting
 209 NTK and $\hat{\Theta}_L$ the empirical finite NTK. Entrywise tensor bounds give

$$\text{Var}((\hat{\Theta}_L)_{ab}) \lesssim L^3 \bar{\varepsilon}, \quad |\mathbb{E}(\hat{\Theta}_L)_{ab} - (\Theta_{\infty,L})_{ab}| \lesssim L^3 \bar{\varepsilon}. \quad (24)$$

210 A layer martingale and Matrix Freedman yield the operator form. Matrix Freedman is the matrix-
 211 valued analogue of Freedman’s martingale concentration inequality: it controls the operator norm
 212 of a sum of conditionally mean-zero random matrices using a bound on the largest increment and
 213 on the predictable quadratic variation. The well-spread data assumption below means that no small
 214 group of examples dominates the random feature directions; equivalently, the matrix variance of the
 215 empirical NTK fluctuation grows like m , not m^2 . This is not true for every deterministic dataset:
 216 duplicated, nearly duplicated, or highly aligned examples can violate it. It is a standard incoherence-
 217 type assumption, and it holds with high probability for many random generic datasets, for example
 218 iid sub-Gaussian or spherical inputs with bounded aspect ratio and no near-collisions, but we keep it
 219 explicit rather than hiding it inside the theorem.

220 **Theorem 6.1** (Finite-network operator concentration). *Assume $L^3 \bar{\varepsilon} \leq c$. Assume further that the*
 221 *data are well-spread in the layerwise feature bases, so the entrywise tensor variance bounds upgrade*
 222 *to a matrix variance proxy of order $mL^3 \bar{\varepsilon}$. With probability at least $1 - \delta$,*

$$\|\hat{\Theta}_L - \Theta_{\infty,L}\|_{\text{op}} \leq CL^{3/2} \sqrt{m \bar{\varepsilon}} \text{polylog}(m, L, 1/\delta) + CmL^3 \bar{\varepsilon} \text{polylog}(m, L, 1/\delta). \quad (25)$$

223 *Without the well-spread data assumption, the same argument still gives a robust bound, but the*
 224 *variance term is weaker:*

$$CL^{3/2} \sqrt{m \bar{\varepsilon}} \text{ is replaced by } CmL^{3/2} \sqrt{\bar{\varepsilon}}.$$

225 *The deterministic bias term $CmL^3 \bar{\varepsilon}$ is unchanged.*

226 *Proof sketch.* Expose the random weights and projectors layer by layer and decompose $\hat{\Theta}_L - \mathbb{E}\hat{\Theta}_L$ as
 227 a matrix martingale. Proposition 5.2 gives conditional entrywise variance $O(L^3 \bar{\varepsilon})$ and conditional bias
 228 $O(L^3 \bar{\varepsilon})$. Under the well-spread data assumption, the predictable quadratic variation is $O(mL^3 \bar{\varepsilon})$ in
 229 operator norm rather than $O(m^2 L^3 \bar{\varepsilon})$. Matrix Freedman gives the first term in (25), with logarithmic
 230 losses from truncation and union over tensor coordinates. The deterministic bias contributes the
 231 second term. Without this assumption, the same argument uses the Frobenius-to-operator bound on
 232 the variance proxy and loses a factor \sqrt{m} . \square

233 If $\mu_{\infty}(L) = \lambda_{\min}(P_m \Theta_{\infty,L} P_m)$ is the limiting centered spectral margin, Weyl’s inequality gives
 234 preservation whenever the right-hand side of (25) is at most $\mu_{\infty}(L)/2$. In the RF-LR bottleneck, the
 235 baseline is not the dense MLP margin; the centered proxy margin can scale as $1/(rL)$. Combining
 236 this with the bottleneck deviation gives the conservative design rule $r \gtrsim mL^3$, up to logarithmic
 237 factors, when the RF-LR perturbation scale is $1/r$.

238 6.1 RF-LR memory and rank-depth rules

239 The RF-LR recursion (5) is not a full-rank-compatible perturbation. Its old-kernel memory obeys an
 240 exact product rule:

$$\Theta_{\text{old} \leftarrow k}^{(L)} = \Theta^{(k)} \prod_{j=k+1}^L \frac{\dot{\Sigma}^{(j)}}{r}. \quad (26)$$

241 At the normalized ReLU edge of chaos, $\hat{\Sigma}^{(j)} = 1/2 + O(j^{-1})$, hence

$$|\Theta_{\text{old} \leftarrow k}^{(L)}| \leq C \left(\frac{1}{2r}\right)^{L-k} \left(\frac{L}{k}\right)^C. \quad (27)$$

242 Low rank therefore shortens old NTK memory geometrically. The cost is that the centered determin-
243 istic margin can shrink with r and L .

244 **Proposition 6.2** (Unprojected and centered RF-LR stability rules). *Consider an RF-LR bottleneck*
245 *whose centered limiting margin satisfies $\mu_{\text{RF}}(L, r) \gtrsim (rL)^{-1}$. If the empirical operator deviation is*
246 *controlled by the conservative scale $L^{3/2}\sqrt{m/r} + mL^3/r$, then spectral preservation follows from*

$$r \gtrsim mL^3 \text{ polylog}(m, L, 1/\delta). \quad (28)$$

247 *If the radial mode is projected out and an angularized centered decomposition improves the variance*
248 *scale by one power of L , then the sufficient rule improves to*

$$r \gtrsim mL^2 \text{ polylog}(m, L, 1/\delta). \quad (29)$$

249 *Proof.* Weyl’s inequality requires the operator deviation to be at most a fixed fraction of $\mu_{\text{RF}}(L, r)$.
250 In the unprojected RF-LR model, the radial contribution remains present in the A -sector, whose tensor
251 scale is cubic in depth. The sufficient inequality is therefore dominated by the mL^3/r contribution
252 after absorbing logarithms and constants, giving (28). In the centered/angularized model, the radial
253 sector is removed before concentration and the dominant tensor variance is quadratic in depth.
254 Repeating the same Weyl argument gives (29). The second statement is conditional on this projection
255 and should not be read as the default RF-LR rule. \square

256 **Training-time drift.** dNTK and ddNTK descendants use the same Gaussian-ReLU atoms and
257 Stiefel-projector cumulants as the NTK itself; differentiation changes deterministic leg labels, not the
258 low-rank random vertex set. The full statement is in Appendix C.

259 **Experiments.** The main figures above give the projector, tensor-constant, memory, and slope
260 diagnostics next to the claims they support. Appendix H reports the remaining true finite-network
261 empirical NTKs, the deep cumulant stress test, and exact full-rank matching.

262 The true finite-network NTK plots, log-log cumulant panels, effective slope table, and depth-200 full-
263 rank matching table are reported in Appendix H. They support the same claims: operator deviations
264 follow the predicted normalization, finite-network cumulants exhibit polynomial effective slopes, and
265 the full-rank endpoint matches the dense NTK to numerical precision.

266 7 Conclusion

267 Low rank changes finite-width NTK theory not only by reducing parameters, but by changing the
268 stability threshold for training. The MLP-compatible model isolates full-rank-compatible rank defects,
269 whereas RF-LR bottlenecks suppress old NTK memory by approximately $(2r)^{-(L-k)}$ while also
270 shrinking the centered deterministic gap. Combining these effects gives the main design message:
271 preserving the NTK edge of convexity has a conservative cubic-depth sufficient rule, $r \gtrsim mL^3$, up to
272 logarithmic factors. The experiments do not contradict this theorem; they show that scalar-output
273 contracted observables can have smaller effective exponents, around 1.5–2.1, suggesting a sharper
274 mode-dependent law between $L^{3/2}$ and L^2 . The rank-width diagrammatics developed here connect
275 these criteria to conditional operator concentration and to finite-network experiments.

276 The consequence is that rank should be treated as a stability budget, not only as a compression
277 knob. If the rank is too small, the model may still have fewer parameters, but its empirical NTK
278 can drift far enough from the dense geometry to lose the spectral margin that makes lazy training
279 predictable. The results also explain why full-rank factorization is harmless, why RF-LR bottlenecks
280 behave differently, and why finite-network experiments can show milder effective depth laws than the
281 conservative theorem. In practice, the theory gives a cautious certificate: it identifies when low rank
282 is expected to preserve the optimization geometry, and where sharper mode-dependent predictions
283 are still needed.

284 To our knowledge, this is the first framework that reconciles low-rank feature geometry with an NTK
285 convexity certificate: low rank can select interpretable bottleneck features while still preserving the
286 spectral conditions that make lazy training predictable.

287 **References**

- 288 [1] A. Jacot, F. Gabriel, and C. Hongler. Neural tangent kernel: convergence and generalization in
289 neural networks. *NeurIPS*, 2018.
- 290 [2] J. Lee, L. Xiao, S. Schoenholz, Y. Bahri, J. Sohl-Dickstein, and J. Pennington. Wide neural
291 networks of any depth evolve as linear models under gradient descent. *NeurIPS*, 2019.
- 292 [3] S. Du, J. Lee, H. Li, L. Wang, and X. Zhai. Gradient descent finds global minima of deep neural
293 networks. *ICML*, 2019.
- 294 [4] Z. Allen-Zhu, Y. Li, and Z. Song. A convergence theory for deep learning via over-
295 parameterization. *ICML*, 2019.
- 296 [5] S. Arora, S. Du, W. Hu, Z. Li, R. Salakhutdinov, and R. Wang. On exact computation with an
297 infinitely wide neural net. *NeurIPS*, 2019.
- 298 [6] Q. Nguyen and M. Mondelli. Global convergence of deep networks with one wide layer
299 followed by pyramidal topology. *NeurIPS*, 2020.
- 300 [7] Q. Nguyen, M. Mondelli, and G. Montúfar. Tight bounds on the smallest eigenvalue of the
301 neural tangent kernel for deep ReLU networks. *ICML*, 2021.
- 302 [8] D. Terjék and D. González-Sánchez. MLPs at the EOC: Concentration of the NTK.
303 arXiv:2501.14724, 2025.
- 304 [9] D. Terjék and D. González-Sánchez. MLPs at the EOC: Spectrum of the NTK.
305 arXiv:2501.13225, 2025.
- 306 [10] G. Yang. Wide feedforward or recurrent neural networks of any architecture are Gaussian
307 processes. *NeurIPS*, 2019.
- 308 [11] G. Yang and E. Hu. Tensor programs IV: Feature learning in infinite-width neural networks.
309 *ICML*, 2021.
- 310 [12] G. Yang, E. J. Hu, I. Babuschkin, S. Sidor, X. Liu, D. Farhi, N. Ryder, J. Pachocki, W. Chen,
311 and J. Gao. Tensor programs V: Tuning large neural networks via zero-shot hyperparameter
312 transfer. *NeurIPS*, 2021.
- 313 [13] G. Yang, D. Yu, C. Zhu, and S. Hayou. Tensor programs VI: Feature learning in infinite-depth
314 neural networks. *ICLR*, 2024.
- 315 [14] B. Bordelon and C. Pehlevan. Self-consistent dynamical field theory of kernel evolution in wide
316 neural networks. *NeurIPS*, 2022.
- 317 [15] S. Mei, A. Montanari, and P.-M. Nguyen. A mean field view of the landscape of two-layer
318 neural networks. *PNAS*, 2018.
- 319 [16] L. Chizat and F. Bach. On the global convergence of gradient descent for over-parameterized
320 models using optimal transport. *NeurIPS*, 2018.
- 321 [17] G. M. Rotskoff and E. Vanden-Eijnden. Neural networks as interacting particle systems:
322 Asymptotic convexity of the loss landscape and universal scaling of the approximation error.
323 arXiv:1805.00915, 2018.
- 324 [18] J. Huang and H.-T. Yau. Dynamics of deep neural networks and neural tangent hierarchy. *ICML*,
325 2020.
- 326 [19] E. Dyer and G. Gur-Ari. Asymptotics of wide networks from Feynman diagrams.
327 arXiv:1909.11304, 2019.
- 328 [20] M. Guillen, P. Misof, and J. E. Gerken. Finite-width neural tangent kernels from Feynman
329 diagrams. arXiv:2508.11522, 2025.
- 330 [21] R. Novak, J. Sohl-Dickstein, and S. Schoenholz. Fast finite width neural tangent kernel. *ICML*,
331 2022.

332 [22] U. Jang, J. D. Lee, and E. K. Ryu. LoRA training in the NTK regime has no spurious local
333 minima. *ICML*, 2024.

334 [23] E. Mosig García, A. Agazzi, and D. Trevisan. Quantitative convergence of trained single layer
335 neural networks to Gaussian processes. arXiv:2509.24544, 2025.

336 **NeurIPS Paper Checklist**

337 **1. Claims**

338 Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s
339 contributions and scope?

340 Answer: [\[Yes\]](#)

341 Justification: The abstract and introduction state the paper’s scope as a finite-width and finite-rank
342 NTK analysis, and the assumptions are made explicit in the theorem statements. The empirical
343 MNIST table is presented only as motivation, not as evidence for the main theoretical claims.

344 Guidelines:

- 345 • The answer [\[N/A\]](#) means that the abstract and introduction do not include the claims made
346 in the paper.
- 347 • The abstract and/or introduction should clearly state the claims made, including the contribu-
348 tions made in the paper and important assumptions and limitations. A [\[No\]](#) or [\[N/A\]](#) answer
349 to this question will not be perceived well by the reviewers.
- 350 • The claims made should match theoretical and experimental results, and reflect how much
351 the results can be expected to generalize to other settings.
- 352 • It is fine to include aspirational goals as motivation as long as it is clear that these goals are
353 not attained by the paper.

354 **2. Limitations**

355 Question: Does the paper discuss the limitations of the work performed by the authors?

356 Answer: [\[Yes\]](#)

357 Justification: The conclusion states the narrow scope of the results, including the restriction to
358 signed/isometric low-rank factors and RF-LR bottlenecks, the status of the endpoint constants,
359 the logarithmic losses, and the open dNTK/ddNTK coefficient problem.

360 Guidelines:

- 361 • The answer [\[N/A\]](#) means that the paper has no limitation while the answer [\[No\]](#) means that
362 the paper has limitations, but those are not discussed in the paper.
- 363 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 364 • The paper should point out any strong assumptions and how robust the results are to violations
365 of these assumptions.
- 366 • The authors should reflect on the scope of the claims made.
- 367 • The authors should reflect on the factors that influence the performance of the approach.
- 368 • The authors should discuss the computational efficiency of the proposed algorithms and how
369 they scale with dataset size.
- 370 • If applicable, the authors should discuss possible limitations of their approach to address
371 problems of privacy and fairness.
- 372 • While the authors might fear that complete honesty about limitations might be used by
373 reviewers as grounds for rejection, reviewers will be specifically instructed to not penalize
374 honesty concerning limitations.

375 **3. Theory assumptions and proofs**

376 Question: For each theoretical result, does the paper provide the full set of assumptions and a
377 complete proof?

378 Answer: [\[Yes\]](#)

379 Justification: The theorem, proposition, and lemma statements specify the finite-rank, projector,
380 ReLU transport, well-spread-data, and RF-LR assumptions. Proof sketches are in the main text,
381 with the tensor hierarchy, projector cumulants, descendant arguments, and endpoint calculations
382 expanded in the appendices.

383 Guidelines:

- 384 • The answer [N/A] means that the paper does not include theoretical results.
- 385 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 386 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 387 • The proofs can either appear in the main paper or the supplemental material.
- 388 • Any informal proof provided in the core of the paper should be complemented by formal
389 proofs provided in appendix or supplemental material.
- 390 • Theorems and lemmas that the proof relies upon should be properly referenced.

391 4. Experimental result reproducibility

392 Question: Does the paper fully disclose all the information needed to reproduce the main
393 experimental results of the paper?

394 Answer: [Yes]

395 Justification: Appendix H gives widths, ranks, depths, sample sizes, dimensions, repetitions,
396 normalizations, and metrics for the true finite-network and proxy experiments. Appendix I gives
397 the MNIST architecture, optimizer, learning rate, batch size, epochs, seed, and clipping rule.

398 Guidelines:

- 399 • The answer [N/A] means that the paper does not include experiments.
- 400 • If the paper includes experiments, a [No] answer to this question will not be perceived well
401 by the reviewers.
- 402 • If the contribution is a dataset and/or model, the authors should describe the steps taken to
403 make their results reproducible or verifiable.
- 404 • Reproducibility can be accomplished through code, detailed instructions, model access, or
405 other means appropriate to the research performed.

406 5. Open access to data and code

407 Question: Does the paper provide open access to the data and code, with sufficient instructions to
408 faithfully reproduce the main experimental results, as described in supplemental material?

409 Answer: [Yes]

410 Justification: An anonymized supplemental code bundle contains the Feynman/NTK experiment
411 scripts and the MNIST scripts, together with a README describing the contents. Synthetic data
412 are generated by the scripts, and MNIST is downloaded through the standard dataset loader.

413 Guidelines:

- 414 • The answer [N/A] means that paper does not include experiments requiring code.
- 415 • Please see the NeurIPS code and data submission guidelines for more details.
- 416 • While code release is encouraged, a justified [No] is acceptable.
- 417 • The instructions should contain the exact command and environment needed to reproduce
418 the results.
- 419 • The authors should provide instructions on data access and preparation.
- 420 • The authors should provide scripts to reproduce all experimental results for the new proposed
421 method and baselines.
- 422 • At submission time, anonymized versions should be released if applicable.

423 6. Experimental setting/details

424 Question: Does the paper specify all the training and test details necessary to understand the
425 results?

426 Answer: [Yes]

427 Justification: The appendices specify the finite-network initialization protocol, uniform Stiefel
428 sampling, autograd NTK computation, ranks, widths, depths, sample sizes, seeds or repetition
429 counts, and MNIST optimizer/training settings.

- 430 Guidelines:
- 431 • The answer [N/A] means that the paper does not include experiments.
 - 432 • The experimental setting should be presented in the core of the paper to a level of detail that
 - 433 is necessary to appreciate the results and make sense of them.
 - 434 • The full details can be provided either with the code, in appendix, or as supplemental
 - 435 material.

436 7. Experiment statistical significance

437 Question: Does the paper report error bars suitably and correctly defined or other appropriate

438 information about the statistical significance of the experiments?

439 Answer: [Yes]

440 Justification: The true finite-network and proxy experiments report repeated-initialization sum-

441 maries, including medians and means, with repetition counts stated in Appendix H. The MNIST

442 snapshot is explicitly described as a compact motivating run rather than a statistical claim.

443 Guidelines:

- 444 • The answer [N/A] means that the paper does not include experiments.
- 445 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
- 446 intervals, or statistical significance tests, at least for the experiments that support the main
- 447 claims.
- 448 • The factors of variability that the error bars are capturing should be clearly stated.
- 449 • The method for calculating the error bars should be explained.
- 450 • The assumptions made should be given.
- 451 • It should be clear whether the error bar is the standard deviation or the standard error of the
- 452 mean.
- 453 • For asymmetric distributions, authors should avoid misleading symmetric error bars.
- 454 • If error bars are reported, the text should explain how they were calculated and reference the
- 455 corresponding figures or tables.

456 8. Experiments compute resources

457 Question: For each experiment, does the paper provide sufficient information on the computer

458 resources needed to reproduce the experiments?

459 Answer: [No]

460 Justification: The paper reports the model sizes and repetition counts but does not yet give precise

461 hardware, memory, wall-clock time, or total compute estimates. The experiments are small

462 PyTorch/autograd studies and MNIST runs, but the exact compute resources should be added

463 before submission if required.

464 Guidelines:

- 465 • The answer [N/A] means that the paper does not include experiments.
- 466 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or
- 467 cloud provider, including relevant memory and storage.
- 468 • The paper should provide the amount of compute required for each of the individual experi-
- 469 mental runs as well as estimate the total compute.
- 470 • The paper should disclose whether the full research project required more compute than the
- 471 experiments reported in the paper.

472 9. Code of ethics

473 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS

474 Code of Ethics?

475 Answer: [Yes]

476 Justification: The work is a theoretical and diagnostic study using synthetic data and MNIST,

477 with no human-subject study, private data, deployment, or high-risk released model.

478 Guidelines:

- 479 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.

- 480 • If the authors answer [No], they should explain the special circumstances that require a
481 deviation from the Code of Ethics.
482 • The authors should make sure to preserve anonymity.

483 10. Broader impacts

484 Question: Does the paper discuss both potential positive societal impacts and negative societal
485 impacts of the work performed?

486 Answer: [N/A]

487 Justification: The work is foundational NTK theory and finite-network diagnostics, with no
488 deployed system, user-facing model, human-subject data, or direct application domain. Its
489 plausible impact is indirect, through better understanding of efficient low-rank training.

490 Guidelines:

- 491 • The answer [N/A] means that there is no societal impact of the work performed.
492 • If the authors answer [N/A] or [No], they should explain why their work has no societal
493 impact or why the paper does not address societal impact.
494 • Examples of negative societal impacts include malicious or unintended uses, fairness consid-
495 erations, privacy considerations, and security considerations.
496 • The conference expects that many papers will be foundational research and not tied to
497 particular applications.
498 • The authors should consider possible harms from intended use, incorrect results, and misuse.
499 • If there are negative societal impacts, authors may discuss mitigation strategies.

500 11. Safeguards

501 Question: Does the paper describe safeguards that have been put in place for responsible release
502 of data or models that have a high risk for misuse?

503 Answer: [N/A]

504 Justification: The paper does not release a high-risk dataset or model. The released material is
505 code for synthetic NTK diagnostics and a standard MNIST benchmark.

506 Guidelines:

- 507 • The answer [N/A] means that the paper poses no such risks.
508 • Released models that have a high risk for misuse or dual-use should be released with
509 necessary safeguards.
510 • Datasets scraped from the Internet could pose safety risks.
511 • Many papers do not require safeguards, but authors should make a best faith effort.

512 12. Licenses for existing assets

513 Question: Are the creators or original owners of assets used in the paper properly credited and
514 are the license and terms of use explicitly mentioned and properly respected?

515 Answer: [No]

516 Justification: The paper uses MNIST through standard loaders and does not redistribute the
517 dataset, but the current manuscript does not explicitly list MNIST license terms. This should be
518 added to the final supplemental documentation if strict asset-license disclosure is required.

519 Guidelines:

- 520 • The answer [N/A] means that the paper does not use existing assets.
521 • The authors should cite the original paper that produced the code package or dataset.
522 • The authors should state which version of the asset is used and, if possible, include a URL.
523 • The name of the license should be included for each asset.
524 • For scraped data, copyright and terms of service should be provided.
525 • If assets are released, license, copyright information, and terms of use should be provided.
526 • For existing datasets that are re-packaged, both the original and derived licenses should be
527 provided.
528 • If this information is not available online, authors are encouraged to reach out to the asset's
529 creators.

530 **13. New assets**

531 Question: Are new assets introduced in the paper well documented and is the documentation
532 provided alongside the assets?

533 Answer: [Yes]

534 Justification: The new supplemental asset is the experiment-code bundle, documented by a
535 README that lists the Feynman/NTK scripts, MNIST scripts, and compact result files.

536 Guidelines:

- 537 • The answer [N/A] means that the paper does not release new assets.
- 538 • Researchers should communicate the details of the dataset/code/model as part of their
539 submissions via structured templates.
- 540 • The paper should discuss whether and how consent was obtained from people whose asset is
541 used.
- 542 • At submission time, remember to anonymize your assets if applicable.

543 **14. Crowdsourcing and research with human subjects**

544 Question: For crowdsourcing experiments and research with human subjects, does the paper
545 include the full text of instructions given to participants and screenshots, if applicable, as well as
546 details about compensation?

547 Answer: [N/A]

548 Justification: The paper does not involve crowdsourcing or human-subject research.

549 Guidelines:

- 550 • The answer [N/A] means that the paper does not involve crowdsourcing nor research with
551 human subjects.
- 552 • Including this information in the supplemental material is fine.
- 553 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or
554 other labor should be paid at least the minimum wage in the country of the data collector.

555 **15. Institutional review board (IRB) approvals or equivalent for research with human subjects**

556 Question: Does the paper describe potential risks incurred by study participants, whether such
557 risks were disclosed to the subjects, and whether Institutional Review Board approvals or equiva-
558 lent approvals were obtained?

559 Answer: [N/A]

560 Justification: The paper does not involve crowdsourcing or human-subject research, so IRB
561 approval is not applicable.

562 Guidelines:

- 563 • The answer [N/A] means that the paper does not involve crowdsourcing nor research with
564 human subjects.
- 565 • Depending on the country in which research is conducted, IRB approval may be required for
566 human-subject research.
- 567 • Procedures may vary significantly between institutions and locations.
- 568 • For initial submissions, do not include any information that would break anonymity.

569 **16. Declaration of LLM usage**

570 Question: Does the paper describe the usage of LLMs if it is an important, original, or non-
571 standard component of the core methods in this research?

572 Answer: [N/A]

573 Justification: LLMs are not an important, original, or non-standard component of the mathematical
574 methods, experiments, or scientific conclusions.

575 Guidelines:

- 576 • The answer [N/A] means that the core method development in this research does not involve
577 LLMs as any important, original, or non-standard components.
- 578 • Please refer to the NeurIPS LLM policy for what should or should not be described.

579 **A Proof Details for Transport and Tensor Hierarchy**

580 This appendix gives the full proof behind Lemma 5.1 and Proposition 5.2. The main text keeps only
581 the proof idea.

582 **Proof of Lemma 5.1.** Write

$$A_\ell := I - \frac{M}{\ell} + E_\ell, \quad \Phi_{L,k} := A_{L-1}A_{L-2}\cdots A_k \quad (L > k),$$

583 and $\Phi_{k,k} := I$. Since $T_2 = 0$, iteration of (14) gives the exact variation-of-constants identity

$$T_L = \sum_{k=2}^{L-1} \Phi_{L,k+1} k^\beta (s + e_k). \quad (30)$$

584 Let $\Phi_{L,k}^0$ denote the unperturbed product with $E_\ell = 0$. Standard finite-dimensional product estimates
585 for regularly varying matrix products give

$$\Phi_{L,k+1}^0 = \exp\left(-M \sum_{j=k+1}^{L-1} \log\left(1 + \frac{1}{j}\right)\right) (I + O(k^{-1})) = L^{-M} k^M + O(L^{-M} k^{M-1}), \quad (31)$$

586 uniformly for $2 \leq k < L$, with the usual harmless logarithmic enlargement when M is not diagonal-
587 izable. The perturbation is summable:

$$\sum_{\ell \geq 2} \|E_\ell\| \leq C \sum_{\ell \geq 2} \frac{\log^p \ell}{\ell^2} < \infty.$$

588 Discrete Duhamel expansion of the product,

$$\Phi_{L,k+1} - \Phi_{L,k+1}^0 = \sum_{j=k+1}^{L-1} \Phi_{L,j+1} E_j \Phi_{j,k+1}^0,$$

589 together with (31) and the summability of E_j , yields

$$\Phi_{L,k+1} = L^{-M} k^M + R_{L,k}, \quad \sum_{k=2}^{L-1} k^\beta \|R_{L,k}\| = O(L^\beta \log^{p+2} L). \quad (32)$$

590 Substituting (32) into (30) gives

$$T_L = L^{-M} \sum_{k=2}^{L-1} k^{M+\beta I} s + O(L^\beta \log^{p+2} L) + L^{-M} \sum_{k=2}^{L-1} k^{M+\beta I} e_k + O(L^\beta \log^{p+q+2} L). \quad (33)$$

591 The error from e_k is $O(L^\beta \log^{q+1} L)$, hence is absorbed by the displayed remainder. The spectral
592 assumption $\Re\lambda(M) > -(\beta + 1)$ makes the matrix Riemann integral convergent. Euler-Maclaurin
593 summation for matrix powers gives

$$L^{-M} \sum_{k=2}^{L-1} k^{M+\beta I} = L^{\beta+1} \int_0^1 u^{M+\beta I} du + O(L^\beta \log L). \quad (34)$$

594 Since

$$\int_0^1 u^{M+\beta I} du = (M + (\beta + 1)I)^{-1},$$

595 we obtain (15). This proves the lemma.

596 **Proof of Proposition 5.2.** The diagrammatic recursion closes on the three blocks

$$\mathcal{V}_\ell, \quad \mathcal{C}_\ell = (D_\ell, F_\ell), \quad \mathcal{A}_\ell = (A_\ell, B_\ell).$$

597 The ReLU transport assumptions imply deterministic block recursions of the form

$$\mathcal{V}_{\ell+1} = \left(I - \frac{M_V}{\ell} + E_\ell^V \right) \mathcal{V}_\ell + \Xi_\ell^V, \quad (35)$$

$$\mathcal{C}_{\ell+1} = \left(I - \frac{M_C}{\ell} + E_\ell^C \right) \mathcal{C}_\ell + B_\ell^C \mathcal{V}_\ell + \Xi_\ell^C, \quad (36)$$

$$\mathcal{A}_{\ell+1} = \left(I - \frac{M_A}{\ell} + E_\ell^A \right) \mathcal{A}_\ell + B_\ell^A \mathcal{C}_\ell + Q_\ell^A[\mathcal{V}_\ell, \mathcal{V}_\ell] + \Xi_\ell^A, \quad (37)$$

598 where each E_ℓ^* satisfies the summable bound in Lemma 5.1. The source tensors $\Xi_\ell^V, \Xi_\ell^C, \Xi_\ell^A$ are
 599 connected rank-state or projector vertices. By Lemma E.1 and Lemma 4.1, their second-order size is
 600 bounded by $C\bar{\varepsilon}$ uniformly in ℓ :

$$\|\Xi_\ell^V\| + \|\Xi_\ell^C\| + \|\Xi_\ell^A\| \leq C\bar{\varepsilon}. \quad (38)$$

601 The deterministic feed-through maps are uniformly bounded in the pairing basis, and the off-diagonal
 602 NTK transport contributes one power of ℓ per transported NTK leg after summation. Equivalently,
 603 the source exponents entering Lemma 5.1 are $\beta = 0$ for \mathcal{V} , $\beta = 1$ for \mathcal{C} , and $\beta = 2$ for \mathcal{A} .

604 For \mathcal{V} , (35) and (38) give Lemma 5.1 with $\beta = 0$, hence

$$\|\mathcal{V}_L\| \leq C\bar{\varepsilon}L.$$

605 For \mathcal{C} , the effective source in (36) is $B_\ell^C \mathcal{V}_\ell + \Xi_\ell^C$. The bound on \mathcal{V}_ℓ gives

$$\|B_\ell^C \mathcal{V}_\ell + \Xi_\ell^C\| \leq C\bar{\varepsilon}\ell,$$

606 so Lemma 5.1 with $\beta = 1$ yields

$$\|\mathcal{C}_L\| \leq C\bar{\varepsilon}L^2.$$

607 For \mathcal{A} , the effective source in (37) is

$$B_\ell^A \mathcal{C}_\ell + Q_\ell^A[\mathcal{V}_\ell, \mathcal{V}_\ell] + \Xi_\ell^A.$$

608 Using the previous two bounds,

$$\|B_\ell^A \mathcal{C}_\ell\| \leq C\bar{\varepsilon}\ell^2, \quad \|Q_\ell^A[\mathcal{V}_\ell, \mathcal{V}_\ell]\| \leq C\bar{\varepsilon}^2\ell^2.$$

609 In the perturbative regime $\bar{\varepsilon} \leq 1$, the quadratic term is bounded by $C\bar{\varepsilon}\ell^2$, and $\Xi_\ell^A \leq C\bar{\varepsilon}$ is smaller.
 610 Lemma 5.1 with $\beta = 2$ gives

$$\|\mathcal{A}_L\| \leq C\bar{\varepsilon}L^3.$$

611 Projecting these vector-block estimates onto the coordinates V, D, F, A, B gives (17). Since the
 612 perturbative NTK mean correction is a finite linear combination of these blocks, its largest contribution
 613 is controlled by $C\bar{\varepsilon}L^3$, which is small when $L^3\bar{\varepsilon} \ll 1$.

614 B Three-Layer Finite-Rank Calculation

615 The simplest place where the low-rank simplification is visible is a three-layer network, i.e. two
 616 ReLU feature maps. Let

$$Q_x = \frac{1}{r} \sum_{a=1}^r X_a^2, \quad Q_y = \frac{1}{r} \sum_{a=1}^r Y_a^2, \quad S = \frac{1}{r} \sum_{a=1}^r X_a Y_a, \quad C_r = \frac{S}{\sqrt{Q_x Q_y}}.$$

617 For the right-factor-only low-rank model, the empirical second hidden-layer kernel has the schematic
 618 form

$$K_{2,r} = \sqrt{Q_x Q_y} \kappa(C_r), \quad \dot{K}_{2,r} = \dot{\kappa}(C_r), \quad \hat{\Theta}_r^{(3)} = K_{2,r} + S \dot{K}_{2,r},$$

619 where κ is the normalized ReLU covariance map.

620 **Proposition B.1** (Delta-method correction for three layers). *For fixed correlation $\rho \in (-1, 1)$,*
 621 *there is an explicit smooth function $\Delta^{(3)}(\rho)$, depending only on $\kappa, \dot{\kappa}$ and the covariance matrix of*
 622 *(X^2, Y^2, XY) , such that*

$$\mathbb{E} \widehat{\Theta}_r^{(3)}(\rho) = \Theta_\infty^{(3)}(\rho) + \frac{1}{r} \Delta^{(3)}(\rho) + O(r^{-2}), \quad \widehat{\Theta}_r^{(3)} - \mathbb{E} \widehat{\Theta}_r^{(3)} = O_{\mathbb{P}}(r^{-1/2}). \quad (39)$$

623 *For scale-invariant ReLU on the diagonal, $\Delta^{(3)}(1) = 0$.*

624 *Proof.* Write $\widehat{\Theta}_r^{(3)} = F(Q_x, Q_y, S)$. By the multivariate central limit theorem and Lemma E.1,

$$(Q_x, Q_y, S) = (1, 1, \rho) + r^{-1/2} \xi + O_{\mathbb{P}}(r^{-1})$$

625 with ξ asymptotically Gaussian and covariance determined by the moments of (X^2, Y^2, XY) . Taylor
 626 expansion gives

$$\mathbb{E} F(Q_x, Q_y, S) = F(1, 1, \rho) + \frac{1}{2r} \text{tr}(\nabla^2 F(1, 1, \rho) \Sigma_\rho) + O(r^{-2}),$$

627 which is (39). The $O_{\mathbb{P}}(r^{-1/2})$ fluctuation is the first-order delta-method term. On the diagonal
 628 $X = Y$, ReLU scale invariance makes the normalized correlation and derivative atoms deterministic
 629 along the radial direction, cancelling the first mean correction. \square

630 C NTK Descendants and Training-Time Drift

631 This appendix records the short argument behind the training-time drift statement. Let $J_c = \nabla_{\theta} f(x_c)$,
 632 $\Xi_{ab;c} = D\Theta_{ab}[J_c]$, and $\Psi_{ab;cd} = D^2\Theta_{ab}[J_c, J_d]$. These descendants add derivative legs to the same
 633 layer maps; they do not introduce new independent random variables.

634 **Claim.** The joint cumulants of (K, Θ, Ξ, Ψ) are generated by the same Gaussian-ReLU atoms and
 635 Stiefel-projector cumulants of $G = (n/r)UU^\top$. Indeed, each derivative with respect to a trainable
 636 right factor $B^{(\ell)}$ inserts a deterministic backpropagated leg and the same frozen projector UU^\top .
 637 Higher derivatives only differentiate ReLU gates and covariance maps, replacing $\kappa, \dot{\kappa}$ by tensors such
 638 as $\ddot{\kappa}$. Since U is fixed, there is no moving-projector derivative. Thus the connected diagrams for
 639 K, Θ, Ξ, Ψ have the same random vertices as the NTK diagrams and differ only in deterministic leg
 640 labels.

641 Under square-loss gradient flow, $\dot{\Theta}_{ab,t} = -\sum_c \Xi_{ab;c,t} e_{c,t}$. The descendant concentration bound
 642 therefore yields, in a lazy window,

$$\|\Theta_t - \Theta_0\|_{\text{op}} \leq Ct \|e_0\|_2 \left(L^3 \sqrt{m\bar{\varepsilon}} + mL^3 \bar{\varepsilon} \right) \text{polylog}(m, L, 1/\delta).$$

643 D Full Proof of the Frozen-Left Decomposition

644 We prove Proposition 3.2. Consider one layer

$$W = \alpha U B, \quad U^\top U = I_r, \quad B \in \mathbb{R}^{r \times d},$$

645 and let $H = \partial \mathcal{L} / \partial W$ be the Euclidean loss gradient in dense-weight coordinates. The Frobenius
 646 inner product is used throughout.

647 First freeze U . For a variation dB ,

$$dW = \alpha U dB, \quad d\mathcal{L} = \langle H, \alpha U dB \rangle = \langle \alpha U^\top H, dB \rangle.$$

648 Hence $\nabla_B \mathcal{L} = \alpha U^\top H$. Gradient descent in B gives

$$\dot{B} = -\eta \alpha U^\top H, \quad \dot{W}_B = \alpha U \dot{B} = -\eta \alpha^2 U U^\top H,$$

649 which is (3). Therefore the B -only tangent directions are exactly the fixed linear subspace

$$\mathcal{T}_B W = \{ \alpha U \Delta B : \Delta B \in \mathbb{R}^{r \times d} \}.$$

650 The corresponding contribution to the empirical NTK is a fixed-projector block, because every dense
651 gradient is projected through UU^\top .

652 Now allow U to vary on the Stiefel manifold $\text{St}(n, r) = \{U : U^\top U = I_r\}$. Its tangent space is

$$T_U \text{St}(n, r) = \{\Delta U : U^\top \Delta U + \Delta U^\top U = 0\}.$$

653 Every tangent vector decomposes as

$$\Delta U = U\Omega + U_\perp K, \quad \Omega^\top = -\Omega, \quad U^\top U_\perp = 0,$$

654 where $U\Omega$ rotates the basis inside $\text{span}(U)$ and $U_\perp K = (I - UU^\top)\Delta U$ moves the subspace itself.
655 For a variation dU ,

$$dW = \alpha dU B, \quad d\mathcal{L} = \langle H, \alpha dU B \rangle = \langle \alpha H B^\top, dU \rangle.$$

656 Thus the Euclidean gradient in U -coordinates is $\alpha H B^\top$. Its normal subspace-moving component is

$$(I - UU^\top)\alpha H B^\top.$$

657 The skew component $U\Omega$ only rotates coordinates inside the same represented subspace: by replacing
658 U with UR and B with $R^\top B$, the product UB is unchanged. This is why the main text states (4) up
659 to rotations inside $\text{span}(U)$ that can be absorbed into B .

660 Keeping only the subspace-moving component, Stiefel gradient descent contributes

$$\dot{U}_\perp = -\eta \alpha (I - UU^\top) H B^\top.$$

661 Mapping this velocity back to dense-weight coordinates gives

$$\dot{W}_U = \alpha \dot{U}_\perp B = -\eta \alpha^2 (I - UU^\top) H B^\top B,$$

662 which is (4). This term is absent when U is frozen. It depends on the current right factor through
663 $B^\top B$, so it cannot be represented by cumulants of the fixed Stiefel projector $G = (n/r)UU^\top$ alone.

664 Finally, if both U and B are trained, the tangent space contains the sum of B -directions $\alpha U \Delta B$,
665 Stiefel subspace directions $\alpha \Delta U B$, and their cross terms in the NTK Gram matrix. Along gradient
666 flow these directions also vary in time, producing dNTK terms involving \dot{U} , \dot{B} , and products such as
667 $B^\top B$. The frozen-left model deliberately removes these moving-subspace interactions, leaving a
668 closed fixed-projector diagrammatic calculus.

669 E Rank Cumulant Power Counting

670 **Lemma E.1** (Rank cumulant power counting). *Fix a layer ℓ and condition on all previous layers.*
671 *For rank channel a , let*

$$Z_a^{(\ell)} := \left\{ (g_a^{(\ell)}(x_u), \sigma(g_a^{(\ell)}(x_u)), \sigma'(g_a^{(\ell)}(x_u)))_{u=1}^m, \text{ and the corresponding derivative atoms} \right\}.$$

672 Here $g_a^{(\ell)}(x_u)$ is the scalar preactivation of rank channel a on sample x_u , after conditioning on layer
673 $\ell - 1$. Thus $Z_a^{(\ell)}$ is the local random object from which one-channel quantities such as σ , σ' , $\sigma'\sigma'$,
674 and derivative insertions are evaluated. The variables $Z_1^{(\ell)}, \dots, Z_r^{(\ell)}$ are conditionally iid. Write
675 them as Z_a , and assume the local functions below have uniformly bounded moments of all orders.
676 For

$$X_i^{(r)} = r^{-a_i} \sum_{a=1}^r f_i(Z_a), \quad i = 1, \dots, s,$$

677 the connected cumulant satisfies

$$\text{cum} \left(X_1^{(r)}, \dots, X_s^{(r)} \right) = r^{1 - \sum_i a_i} \text{cum} (f_1(Z), \dots, f_s(Z)). \quad (40)$$

678 In particular, normalized empirical averages have s -point cumulants of order r^{1-s} .

679 *Proof.* By multilinearity,

$$\text{cum}(X_1^{(r)}, \dots, X_s^{(r)}) = r^{-\sum_i a_i} \sum_{a_1, \dots, a_s=1}^r \text{cum}(f_1(Z_{a_1}), \dots, f_s(Z_{a_s})).$$

680 Independence kills every term whose indices are not all in the same connected block. Only $a_1 =$
681 $\dots = a_s$ contributes, giving r identical terms and hence (40). \square

682 F Projector Cumulants Beyond Second Order

683 This appendix records the higher-order cumulant scaling of the same projector G used in the main
 684 proof. For $G = (n/r)UU^\top$, the mean is $\mathbb{E}G = I_n$, the trace is deterministic, and every connected
 685 cumulant containing the trace direction vanishes. Orthogonal Weingarten calculus gives, for fixed
 686 cumulant order s ,

$$\text{cum}(G_{i_1 j_1}, \dots, G_{i_s j_s}) = O_s(r^{1-s}) \quad (41)$$

687 uniformly when $r \leq n$ and $r \rightarrow \infty$, with the full-rank endpoint cancelling after replacing r^{-1} by the
 688 defect scale $\gamma_{n,r}$ at second order. Diagrammatically, this is why a connected rank vertex is counted by
 689 the number of connected projector insertions, not simply by the number of factors U in the algebraic
 690 expression.

691 G Low-Rank Diagrammatic Rules

692 This appendix gives the graphical rules behind the main tensor recursions. This appendix records the
 693 graphical calculus underlying the perturbative statements in the main text. The external objects are the
 694 same as in finite-width NTK Feynman diagrams: NNGP/preactivation legs, NTK fluctuation legs, and
 695 derivative-NTK descendants. The difference is internal. In the low-rank model, stochastic vertices
 696 are empirical rank cumulants and frozen-projector cumulants rather than unconstrained neural-index
 697 contractions.

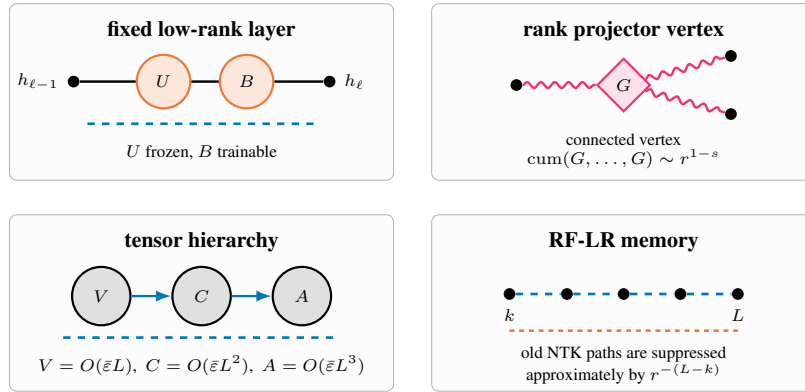


Figure 5: Main diagrammatic objects. Solid black lines denote NNGP/preactivation transport, dashed blue lines denote NTK legs, orange vertices denote low-rank factors, pink diamonds denote Stiefel-projector cumulants, and gray blobs denote transported tensor cumulants.

698 G.1 Dense NTK rules and low-rank replacement

699 This subsection makes explicit how the finite-width NTK Feynman rules are imported and modified
 700 [20]. The dense calculus has the following ingredients.

- 701 1. External filled dots are the observed sample labels. Solid external lines denote preactivation or
 702 NNGP objects $z_\alpha, K_{\alpha\beta}$. Colored dashed external lines denote NTK fluctuations $\widehat{\Delta}\Theta_{\alpha\beta}$. Double
 703 or multi-dashed descendants denote dNTK and ddNTK insertions.
- 704 2. Cubic vertices attach two external legs to one internal line. The decoration of the internal line
 705 is the local Gaussian-ReLU atom generated by Taylor expansion of the layer map, for example
 706 $\widehat{\Delta}\Omega_{i,\alpha\beta}, \sigma_{i,\alpha}\sigma'_{i,\beta}$, or $\sigma'_{i,\alpha}\sigma'_{i,\beta}$. In the dense network each such vertex carries the usual n_ℓ^{-1}
 707 counting from the neural index.
- 708 3. A white propagator is a Gaussian expectation $\langle \cdot \rangle_{K^{(\ell)}}$. It contracts the decorations of internal
 709 lines under the infinite-width NNGP covariance. The selection rules say that propagators connect
 710 allowed internal lines, equal neural indices are identified, undecorated dashed NTK lines do
 711 not enter the Gaussian expectation, and paired tangent colors contribute the corresponding $\Theta_{\alpha\beta}^{(\ell)}$
 712 factor.

- 713 4. Quartic blobs encode the rank-four tensors that close the order- $1/n$ dense hierarchy. The V
714 blob has four K -legs, D, F have two K -legs and two Θ -legs, and A, B have four Θ -legs. The
715 same external grammar also defines the higher-derivative tensors P, Q, R, S, T, U for dNTK and
716 ddNTK.
- 717 5. To compute a cumulant or mean correction, one draws all connected diagrams with the prescribed
718 external legs and the desired order in $1/n$, applies the propagator selection rules, translates each
719 diagram into a Gaussian expectation of derivatives of local atoms, and sums the admissible
720 diagrams.

721 Our low-rank rules keep this external grammar exactly. The replacement happens only inside the
722 stochastic source:

$$\begin{aligned} \text{dense neural-index source} &\longrightarrow \text{rank-state cumulants } r_\ell^{1-s} \kappa_s^{(\ell)} \\ &\quad \text{and projector cumulants } \text{cum}(G^{(\ell)}, \dots, G^{(\ell)}) \\ \frac{1}{n_\ell} \text{ dense second cumulant} &\longrightarrow \frac{1}{n_\ell} \text{ dense part} + \gamma_{n_\ell, r_\ell} \text{ projector part} + \frac{1}{r_\ell} \text{ rank-state part.} \end{aligned}$$

723 At second order the projector contribution is exactly Lemma 4.1. It is proportional to

$$\gamma_{n,r} = \frac{n(n-r)}{r(n-1)(n+2)},$$

724 so it vanishes at $r = n$. This is the main structural difference from a naive dense replacement $n \mapsto r$:
725 dense Feynman rules have no reason to cancel at full rank, whereas the isometric low-rank projector
726 rules do.

727 G.2 Worked example: one NTK-mean correction

728 We now compute one representative diagram in algebraic form. Let $M_A^{(\ell)} = r_\ell^{-1} \sum_{a=1}^{r_\ell} \phi_A^{(\ell)}(Z_a)$ be
729 a rank-state observable, where A labels a local atom such as $\sigma\sigma, \sigma\sigma'$, or $\sigma'\sigma'$. Suppose the one-layer
730 NTK update can be written locally as a smooth map

$$\widehat{\Theta}_{12}^{(\ell+1)} = \Phi_{12}(M^{(\ell)}, G^{(\ell)}) \text{ transported lower-layer terms.}$$

731 Taylor expand Φ_{12} around the deterministic state (\bar{M}, I) . The connected rank-state part of the mean
732 correction is

$$\mathbb{E} \widehat{\Theta}_{12}^{(\ell+1)} - \Theta_{12}^{(\ell+1)} = \frac{1}{2r_\ell} \sum_{A,B} \partial_A \partial_B \Phi_{12}(\bar{M}, I) \kappa_{AB}^{(\ell)} + \dots, \quad (42)$$

733 where

$$\kappa_{AB}^{(\ell)} = \text{cum}(\phi_A^{(\ell)}(Z), \phi_B^{(\ell)}(Z)).$$

734 This is the low-rank analogue of a Misof propagator diagram: their white Gaussian propagator
735 evaluates a dense neural-index contraction; here the connected empirical-rank vertex evaluates the
736 cumulant of one rank channel and contributes the explicit factor r_ℓ^{-1} .

737 The frozen projector gives a second, full-rank-compatible correction. If $G = I + \Delta G$, then the
738 leading projector diagram is

$$\begin{aligned} &\frac{1}{2} \sum_{ij,kl} \partial_{G_{ij}} \partial_{G_{kl}} \Phi_{12}(\bar{M}, I) \text{Cov}(G_{ij}, G_{kl}) \\ &= \frac{\gamma_{n_\ell, r_\ell}}{2} \sum_{ij,kl} \partial_{G_{ij}} \partial_{G_{kl}} \Phi_{12}(\bar{M}, I) \left(\delta_{ik} \delta_{jl} + \delta_{il} \delta_{jk} - \frac{2}{n_\ell} \delta_{ij} \delta_{kl} \right). \end{aligned} \quad (43)$$

739 Equations (42) and (43) are the same diagrammatic operation as in the dense NTK calculus: choose
740 the external Θ_{12} legs, attach the allowed local vertices, evaluate the connected stochastic source, and
741 sum. The only change is the source value. Dense Misof diagrams count neural-index loops by $1/n$;
742 our diagrams count rank-channel cumulants by $1/r$ and Stiefel-projector cumulants by $\gamma_{n,r}$.

743 For comparison, the dense next-to-leading NTK mean calculation has five admissible diagrams [20]:
 744 transported $\Theta^{\{1\}}$, transported $K^{\{1\}}$, and source diagrams from V, D, F . In our notation this becomes
 745 the schematic low-rank expansion

$$\Theta_{\text{LR},12}^{\{1\}(\ell+1)} = \mathsf{T}_{\Theta}[\Theta_{\text{LR}}^{\{1\}(\ell)}]_{12} \mathsf{T}_K[K_{\text{LR}}^{\{1\}(\ell)}]_{12} \mathsf{S}_V[V_{\ell}]_{12} \mathsf{S}_D[D_{\ell}]_{12} \mathsf{S}_F[F_{\ell}]_{12}, \quad (44)$$

746 but the tensors V, D, F now decompose into dense finite-width, rank-state, and projector-defect
 747 sources. Thus the external five-diagram grammar is imported unchanged, while the internal weights
 748 are changed so that exact full-rank matching is preserved.

749 G.3 Summary of the low-rank results encoded by the rules

750 The diagrammatic rules above encode the paper’s main results as follows.

- 751 1. Exact full-rank matching: when $r_{\ell} = n_{\ell}$, $U^{(\ell)}$ is orthogonal, $G = I$, every centered projector
 752 vertex vanishes, and the factorized NTK equals the dense NTK pathwise.
- 753 2. Stiefel projector covariance: the second projector vertex is $\text{Cov}(G_{ij}, G_{kl}) = \gamma_{n,r}(\delta_{ik}\delta_{jl} +$
 754 $\delta_{il}\delta_{jk} - 2\delta_{ij}\delta_{kl}/n)$, with $\gamma_{n,n} = 0$.
- 755 3. Rank cumulant power counting: a connected s -point rank-state vertex contributes r^{1-s} . Normal-
 756 ized rank averages therefore have s -point cumulants of order r^{1-s} .
- 757 4. Tensor hierarchy: the imported V, C, A external grammar remains triangular, with $V = O(\varepsilon L)$,
 758 $C = O(\varepsilon L^2)$, and $A = O(\varepsilon L^3)$ under the stated ReLU transport assumptions.
- 759 5. Endpoint constants: after pairing-basis extraction, the isolated off-diagonal constants are 5 for V ,
 760 $21/20$ for the mixed C sector, and $173/720$ for the A sector.
- 761 6. Operator criterion: the tensor bounds imply $\|\widehat{\Theta}_L - \Theta_{\infty,L}\|_{\text{op}} \lesssim L^{3/2}\sqrt{m\varepsilon} + mL^3\varepsilon$, up to
 762 logarithmic factors and the well-spread data condition.
- 763 7. RF-LR memory and rank-depth rule: in the bottleneck model, old NTK lines are multiplied by
 764 $\prod_{j=k+1}^L \dot{\Sigma}^{(j)}/r$, giving geometric memory suppression and the conservative stability rule $r \gtrsim$
 765 mL^3 , with the centered/angularized mL^2 rule only under the additional projection hypothesis.

766 G.4 Feynman diagram building blocks and tensor assembly

767 The diagrams below should be read as Feynman diagrams, not as decorative notation. A filled dot
 768 is an observed sample label. A solid external leg represents a preactivation or NNGP entry K_{ab} . A
 769 dashed colored leg represents an NTK fluctuation $\Delta\Theta_{ab}$. Double or multi-dashed legs represent
 770 dNTK and ddNTK descendants. The tensor is determined by the external legs:

$$\begin{aligned} K, K &\rightsquigarrow V, \\ K, \Theta &\rightsquigarrow D, F, \\ \Theta, \Theta &\rightsquigarrow A, B, \\ \text{d}\Theta, K &\rightsquigarrow P, Q, \\ \text{dd}\Theta \text{ external legs} &\rightsquigarrow R, S, T, U. \end{aligned}$$

771 The ordering of sample labels chooses the pairing coordinate. For example, the two K -leg cumulants
 772 $\text{cum}(\delta K_{12}, \delta K_{34})$, $\text{cum}(\delta K_{13}, \delta K_{24})$, and $\text{cum}(\delta K_{14}, \delta K_{23})$ give V_{1234} , V_{1324} , and V_{1423} . Simi-
 773 larly, one solid pair and one dashed pair give D_{1234} , F_{1324} , F_{1423} , while two dashed pairs give A_{1234} ,
 774 B_{1324} , B_{1423} .

775 Once the external legs are fixed, every connected diagram is built from the same three operations:
 776 transport an older tensor through deterministic K - or Θ -leg derivatives, feed a lower block into a
 777 higher block, or attach a fresh low-rank source. The fresh source is where our rules differ from the
 778 dense Misof rules: dense neural-index contractions are replaced by rank-state cumulants $r^{1-s}\kappa_s$ and
 779 Stiefel-projector cumulants, whose second-order scale is $\gamma_{n,r}$.

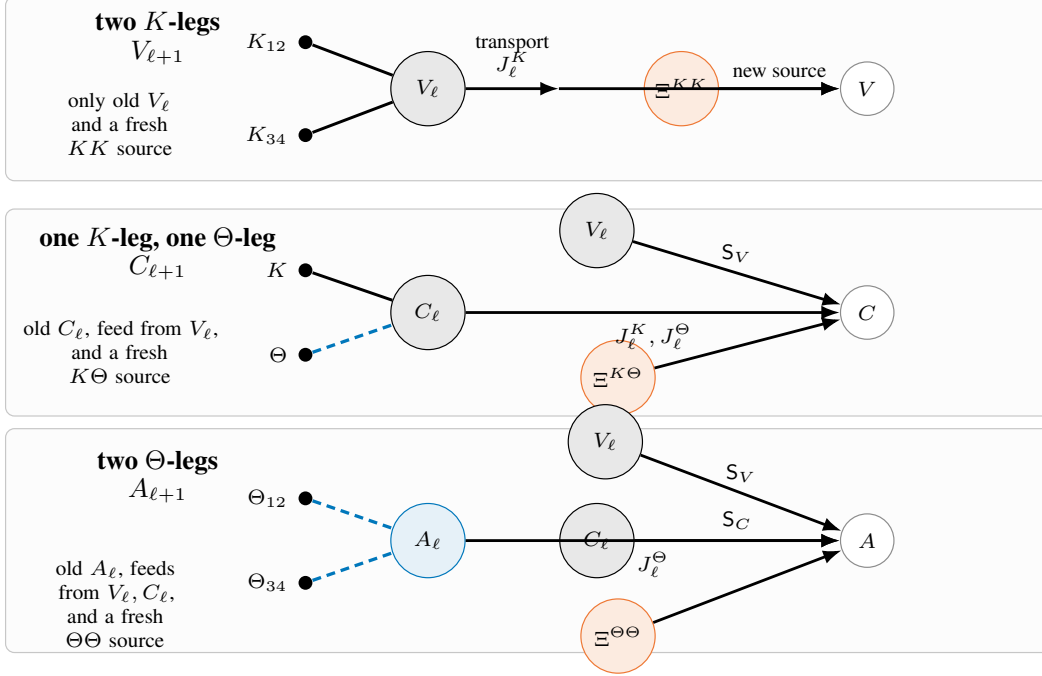


Figure 6: Diagrammatic reduction of the tensor hierarchy. External K - and Θ -legs choose the block being computed. Connected terms either transport an older tensor, feed a lower block into a higher one, or create a fresh low-rank source Ξ .

780 For the NTK-NTK block this gives the schematic reduction

$$\begin{aligned}
\text{cum}(\delta\Theta_{12}^{(\ell+1)}, \delta\Theta_{34}^{(\ell+1)}) &= \underbrace{J_{\ell}^{\Theta} \text{cum}(\delta\Theta_{12}^{(\ell)}, \delta\Theta_{34}^{(\ell)})(J_{\ell}^{\Theta})^{\top}}_{\text{transported old NTK-NTK diagram}} \\
&+ \underbrace{S_C[C_{\ell}] + S_V[V_{\ell}]}_{\text{mixed diagrams feeding into NTK-NTK}} + \underbrace{\Xi_{\ell}^{\Theta\Theta}}_{\text{new connected rank/projector source}}. \tag{45}
\end{aligned}$$

781 The same operation with two K -legs gives $V_{\ell+1}$, and with one K -leg and one Θ -leg gives $C_{\ell+1}$. Thus
782 the apparently large expansion collapses because every connected diagram has only two possible
783 roles: it either transports an older tensor through J_{ℓ}^K or J_{ℓ}^{Θ} , or it creates a fresh source Ξ . This is the
784 mechanism behind the triangular recursion in (11).

785 G.5 Basic tensor recursions

786 The diagrammatic rules produce the same external tensor hierarchy as the dense finite-width expansion,
787 but with low-rank sources. Let J_{ℓ}^K and J_{ℓ}^{Θ} denote the deterministic linear transports for NNGP
788 and NTK legs. The NNGP covariance block satisfies schematically

$$V_{\ell+1} = J_{\ell}^K V_{\ell} (J_{\ell}^K)^{\top} + \Xi_{\ell}^{KK}, \quad \Xi_{\ell}^{KK} = r_{\ell}^{-1} \kappa_2(S^K, S^K) + \gamma_{n_{\ell}, r_{\ell}} \mathfrak{h}_{\ell}^{KK}. \tag{46}$$

789 For the mixed NNGP-NTK block $C = (D, F)$,

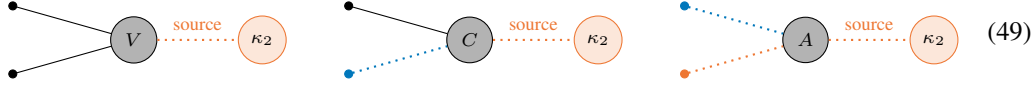
$$C_{\ell+1} = J_{\ell}^K C_{\ell} (J_{\ell}^{\Theta})^{\top} + S_V[V_{\ell}] + \Xi_{\ell}^{K\Theta}. \tag{47}$$

790 For the NTK-NTK block $A = (A, B)$,

$$A_{\ell+1} = J_{\ell}^{\Theta} A_{\ell} (J_{\ell}^{\Theta})^{\top} + S_C[C_{\ell}] + S_V[V_{\ell}] + \Xi_{\ell}^{\Theta\Theta}. \tag{48}$$

791 The sources Ξ^{KK} , $\Xi^{K\Theta}$, and $\Xi^{\Theta\Theta}$ are precisely the rank and projector cumulant vertices with the
792 corresponding external legs. The next display should be read only as a source dictionary. It is not
793 a separate calculation: it says which connected second cumulant is injected into each tensor block.

794 The rank/projector blob labelled κ_2 contributes the small factor r^{-1} , $1/n$, or $\gamma_{n,r}$. The number of
 795 external Θ -legs then determines how many NTK transports are accumulated over depth, and therefore
 796 whether the block grows like L , L^2 , or L^3 .



797 Equivalently, the diagrammatic message is

External legs	Tensor block	Fresh source size	Depth accumulation
K, K	$V = \text{cum}(K, K)$	$O(\bar{\varepsilon})$	no transported NTK leg, hence $O(\bar{\varepsilon}L)$
K, Θ	$C = (D, F)$	$O(\bar{\varepsilon})$	one transported NTK leg, hence $O(\bar{\varepsilon}L^2)$
Θ, Θ	$A = (A, B)$	$O(\bar{\varepsilon})$	two transported NTK legs, hence $O(\bar{\varepsilon}L^3)$

799 Thus the diagrams separate two effects that are easy to confuse. The orange/rank blob gives the
 800 perturbative smallness; it does not decide the power of L . The external K - and Θ -legs decide the
 801 depth power through deterministic transport.

802 G.6 First NTK mean correction

803 The familiar finite-width correction to the NTK mean has the same external five-term form:

$$\Theta_{\ell+1}^{\{1\}} = \mathsf{T}_{\Theta}[\Theta_{\ell}^{\{1\}}] + \mathsf{T}_K[K_{\ell}^{\{1\}}] + \mathsf{S}_V[V_{\ell}] + \mathsf{S}_D[D_{\ell}] + \mathsf{S}_F[F_{\ell}]. \quad (50)$$

804 The low-rank difference is that the source tensors V, D, F are further resolved into the rank-state and
 805 projector vertices above. Thus the low-rank calculus does not change the external NTK grammar; it
 806 changes the internal expansion parameter from neural-index contractions to empirical-rank cumulants
 807 and frozen-projector cumulants.

808 G.7 Full tensor dictionary

809 For four sample labels 1, 2, 3, 4, write

$$\delta K_{ab}^{(\ell)} := K_{ab}^{(\ell)} - \mathbb{E}K_{ab}^{(\ell)}, \quad \delta \Theta_{ab}^{(\ell)} := \Theta_{ab}^{(\ell)} - \mathbb{E}\Theta_{ab}^{(\ell)}.$$

810 All coordinates below are centered two-point cumulants of these pair observables. Since the variables
 811 are centered, $\text{cum}(X, Y) = \mathbb{E}[XY]$. In particular, $V_{1234}^{(\ell)} = \mathbb{E}[\delta K_{12}^{(\ell)} \delta K_{34}^{(\ell)}]$, $D_{1234}^{(\ell)} = \mathbb{E}[\delta K_{12}^{(\ell)} \delta \Theta_{34}^{(\ell)}]$,
 812 $F_{1324}^{(\ell)} = \mathbb{E}[\delta K_{13}^{(\ell)} \delta \Theta_{24}^{(\ell)}]$, $A_{1234}^{(\ell)} = \mathbb{E}[\delta \Theta_{12}^{(\ell)} \delta \Theta_{34}^{(\ell)}]$, and $B_{1324}^{(\ell)} = \mathbb{E}[\delta \Theta_{13}^{(\ell)} \delta \Theta_{24}^{(\ell)}]$. The rank-four tensor
 813 basis is organized by the three pairings (12|34), (13|24), and (14|23). The NNGP-NNGP block is

$$\mathcal{V}_{\ell} = \begin{pmatrix} V_{1234}^{(\ell)} \\ V_{1324}^{(\ell)} \\ V_{1423}^{(\ell)} \end{pmatrix} := \begin{pmatrix} \text{cum}(\delta K_{12}^{(\ell)}, \delta K_{34}^{(\ell)}) \\ \text{cum}(\delta K_{13}^{(\ell)}, \delta K_{24}^{(\ell)}) \\ \text{cum}(\delta K_{14}^{(\ell)}, \delta K_{23}^{(\ell)}) \end{pmatrix}. \quad (51)$$

814 The mixed NNGP-NTK block is

$$\mathcal{C}_{\ell} = \begin{pmatrix} D_{1234}^{(\ell)} \\ F_{1324}^{(\ell)} \\ F_{1423}^{(\ell)} \end{pmatrix} := \begin{pmatrix} \text{cum}(\delta K_{12}^{(\ell)}, \delta \Theta_{34}^{(\ell)}) \\ \text{cum}(\delta K_{13}^{(\ell)}, \delta \Theta_{24}^{(\ell)}) \\ \text{cum}(\delta K_{14}^{(\ell)}, \delta \Theta_{23}^{(\ell)}) \end{pmatrix}. \quad (52)$$

815 The NTK-NTK block is

$$\mathcal{A}_{\ell} = \begin{pmatrix} A_{1234}^{(\ell)} \\ B_{1324}^{(\ell)} \\ B_{1423}^{(\ell)} \end{pmatrix} := \begin{pmatrix} \text{cum}(\delta \Theta_{12}^{(\ell)}, \delta \Theta_{34}^{(\ell)}) \\ \text{cum}(\delta \Theta_{13}^{(\ell)}, \delta \Theta_{24}^{(\ell)}) \\ \text{cum}(\delta \Theta_{14}^{(\ell)}, \delta \Theta_{23}^{(\ell)}) \end{pmatrix}. \quad (53)$$

816 This is the tensor dictionary behind the compact V, C, A notation in the main text: $C = (D, F)$ and
 817 $A = (A, B)$. The dense finite-width Feynman sectors remain the same, but their low-rank stochastic
 818 vertices are cumulants of the frozen projector $G = (n/r)UU^{\top}$.

819 **G.8 Pairing extraction**

820 A rank-four orthogonally invariant tensor decomposes as

$$C_{i_1 i_2 i_3 i_4} = c_1 \delta_{i_1 i_2} \delta_{i_3 i_4} + c_2 \delta_{i_1 i_3} \delta_{i_2 i_4} + c_3 \delta_{i_1 i_4} \delta_{i_2 i_3}. \quad (54)$$

821 Let B_1, B_2, B_3 be the three pairings. The normalized pairing Gram matrix is

$$\bar{G}_N = \begin{pmatrix} 1 & 1/N & 1/N \\ 1/N & 1 & 1/N \\ 1/N & 1/N & 1 \end{pmatrix}, \quad \bar{G}_N^{-1} = \frac{N}{N-1} I_3 - \frac{N}{(N-1)(N+2)} \mathbf{1}\mathbf{1}^\top. \quad (55)$$

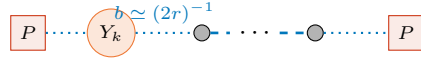
822 If $h = (\langle C, B_1 \rangle, \langle C, B_2 \rangle, \langle C, B_3 \rangle)$ are measured contractions, the isolated tensor coefficients are
 823 $(c_1, c_2, c_3)^\top = \bar{G}_N^{-1} h$. This is why direct scalar-output cumulants should not be identified with the
 824 endpoint coefficients 5, 21/20, and 173/720 without a pairing-basis extraction.

825 **G.9 Centered RF-LR memory lemma**

826 For RF-LR, an old NTK line transported from layer k to L carries

$$B_{L:k} = \prod_{s=k+1}^L \frac{1}{r} \dot{\Sigma}^{(s)}. \quad (56)$$

827 At ReLU edge of chaos, $\dot{\Sigma}^{(s)} \rightarrow 1/2$, so $\|B_{L:k}\| \lesssim (2r)^{-(L-k)}$, up to polynomial endpoint factors.
 828 The corresponding centered transport diagram is



$$\boxed{P} \cdots \overset{b \simeq (2r)^{-1}}{\circlearrowleft} Y_k \cdots \circlearrowleft \cdots \circlearrowleft \cdots \boxed{P} \quad (57)$$

829 Let $P = I - \mathbf{1}\mathbf{1}^\top/m$ and $Z_L = P(\hat{K}_L - K_L)P$. The centered fluctuation admits the linearized
 830 expansion

$$Z_L = \sum_{k=1}^L B_{L:k+1} Y_k + \sum_{k=1}^L B_{L:k+1} R_k, \quad (58)$$

831 where Y_k is a centered fresh rank source and R_k is quadratic in earlier fluctuations.

832 **Lemma G.1** (Centered RF-LR memory). *Assume, conditionally on the past,*

$$\|\mathbb{E}[Y_k^2 \mid \mathcal{F}_{k-1}]\|_{\text{op}} \leq C \frac{m\varepsilon_{\text{RF}}}{r^2}, \quad \|Y_k\|_{\psi_1} \leq C \frac{\sqrt{\varepsilon_{\text{RF}}}}{r}, \quad (59)$$

833 and $\|B_{L:k}\| \leq q^{L-k}$ for some $q < 1$. If $\|R_k\|_{\text{op}} \leq C m\varepsilon_{\text{RF}}/r^2$ with high probability, then, up to
 834 logarithmic factors,

$$\|Z_L\|_{\text{op}} \lesssim \frac{\sqrt{m\varepsilon_{\text{RF}} \log(m/\delta)}}{r} + \frac{m\varepsilon_{\text{RF}} \log(m/\delta)}{r^2} \quad (60)$$

835 with probability at least $1 - \delta$.

836 *Proof.* Set $X_k = B_{L:k+1} Y_k$. The predictable matrix variance satisfies

$$\left\| \sum_{k=1}^L \mathbb{E}[X_k^2 \mid \mathcal{F}_{k-1}] \right\|_{\text{op}} \leq C \frac{m\varepsilon_{\text{RF}}}{r^2} \sum_{k=1}^L q^{2(L-k)} \leq C' \frac{m\varepsilon_{\text{RF}}}{r^2}. \quad (61)$$

837 Matrix Freedman or Bernstein gives the first term in (60). For the quadratic remainder,

$$\left\| \sum_{k=1}^L B_{L:k+1} R_k \right\|_{\text{op}} \leq C \frac{m\varepsilon_{\text{RF}}}{r^2} \sum_{k=1}^L q^{L-k} \leq C' \frac{m\varepsilon_{\text{RF}}}{r^2}. \quad (62)$$

838 \square

839 Combining Lemma G.1 with Weyl's inequality gives the centered/angularized criterion. If the centered
 840 deterministic margin satisfies $\lambda_{\min}(PK_L P) \gtrsim (rL)^{-1}$, and if $\varepsilon_{\text{RF}} \simeq 1/r$, then the fluctuation bound
 841 is below the margin when

$$r \gtrsim mL^2 \quad (63)$$

842 up to logarithms. Without removing the radial constant mode, a diffusive \sqrt{L} factor can reappear,
 843 giving the conservative cubic-depth criterion used for the unprojected RF-LR statement.

844 H Extended True Finite-Network Experiments

845 This appendix keeps the detailed diagnostics out of the main ten-page narrative while preserving the
846 evidence used to interpret the finite-network checks.

847 H.1 Experimental methodology

848 All true finite-network experiments use scalar-output ReLU networks with the signed/isometric
849 factorization $W^{(\ell)} = \sigma_w n_{\ell-1}^{-1/2} \sqrt{n_{\ell}/r_{\ell}} U^{(\ell)} B^{(\ell)}$. The left factors $U^{(\ell)}$ are sampled uniformly from
850 the Stiefel manifold and held fixed; the empirical NTK is computed by PyTorch autograd with respect
851 to the trainable right factors $B^{(\ell)}$ and the readout. Inputs are synthetic, fixed across the compared
852 initializations, and normalized to unit Euclidean norm. The reported quantities are not training
853 accuracies: they are initialization-time NTK and NNGP diagnostics designed to test the rank-width
854 perturbation theory.

855 For the clean true-NTK operator and rank-defect sweep, we use width $n = 128$, sample size $m = 8$,
856 input dimension $d = 16$, ranks $r \in \{8, 16, 32, 64, 128\}$, and depths $L \in \{2, 3, 4, 5, 6, 8, 10\}$. For
857 each pair (r, L) , 200 independent initializations are drawn. For a fixed (r, L) , we compute the
858 empirical NTK Gram matrix $\hat{\Theta}$, subtract its initialization mean across seeds, and report the median
859 and mean operator norm of the fluctuation. The normalization in Figure 7 is $L^{3/2} \sqrt{m\varepsilon}$, with
860 $\varepsilon = 1/n + \gamma_{n,r}$. The rank-defect panel fixes the largest depth in this sweep and varies r , checking
861 that the excess low-rank contribution follows the defect scale and disappears at $r = n$. Exact full-rank
862 matching is tested separately at $n = r = 128$, $m = 8$, $d = 16$, and depths up to $L = 200$; because
863 the theorem is pairwise, this test measures numerical precision rather than a statistical average.

864 For direct finite-network cumulants, we again use $n = 128$, $m = 8$, $d = 16$, and ranks
865 $r \in \{8, 16, 32, 64, 128\}$. For every (r, L) , 2000 independent initializations are used. The stan-
866 dard long sweep uses depths $L \in \{2, 3, 4, 5, 6, 8, 10, 12, 14, 16, 20\}$. The deep stress test uses
867 $L \in \{20, 50, 100, 200\}$. For each initialization we compute the final hidden-feature Gram matrix
868 $K = h_L h_L^{\top}/n$ and the empirical NTK $\Theta = JJ^{\top}$. Across seeds we estimate the centered second
869 cumulants $V = \text{cum}(K, K)$, $C = \text{cum}(K, \Theta)$, and $A = \text{cum}(\Theta, \Theta)$ on fixed same-off-diagonal
870 and disjoint-off-diagonal sample pairings. The log-log slopes in Table 2 are ordinary least-squares
871 fits of $\log |\text{cumulant}|$ against $\log L$ over the listed depth set for each rank. These slopes are diagnos-
872 tic effective exponents of contracted scalar-output observables; they are not estimates of isolated
873 pairing-basis tensor constants.

874 The proxy experiments in the main text are separate Monte Carlo checks of closed-form quantities.
875 The Stiefel-projector covariance experiment uses $n = 64$, ranks $r \in \{4, 8, 16, 32, 64\}$, and 50,000
876 Stiefel samples to compare $\text{Cov}(G_{12}, G_{12})$ with $\gamma_{n,r}$. The scalar endpoint-constant experiment
877 uses $\varepsilon = 10^{-3}$, depths $L \in \{8, 12, 16, 24, 32, 48, 64, 96, 128, 192, 256\}$, and 200,000 Monte Carlo
878 samples to check convergence of the normalized V , C , and A proxies to 5, 21/20, and 173/720. The
879 operator proxy uses $m = 64$ and 300 repetitions per depth to check stability of the normalization
880 $L^{3/2} \sqrt{m\varepsilon}$. The RF-LR memory panel uses ranks 2, 4, 8, 16 and tracks the product of old-NTK
881 transport factors over ages up to 24 layers.

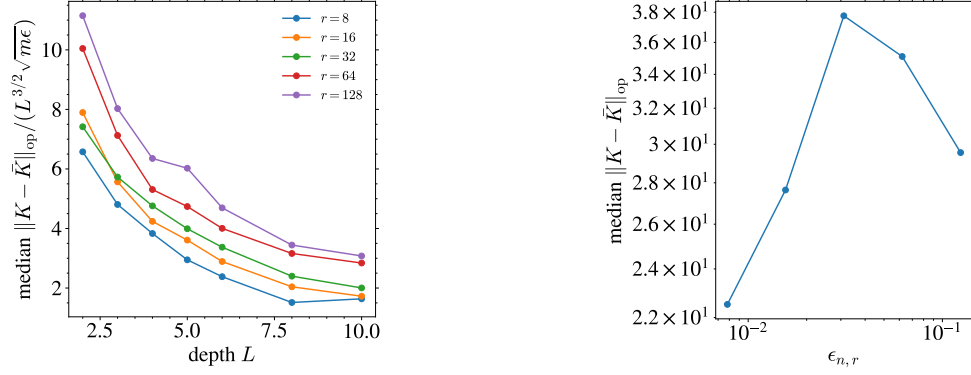


Figure 7: True finite-network NTK checks computed by autodiff. Left: empirical operator deviations follow the predicted normalization. Right: the low-rank excess deviation tracks the rank-defect scale; at the full-rank endpoint the rank-defect contribution disappears, while ordinary finite-width and numerical residuals can remain.

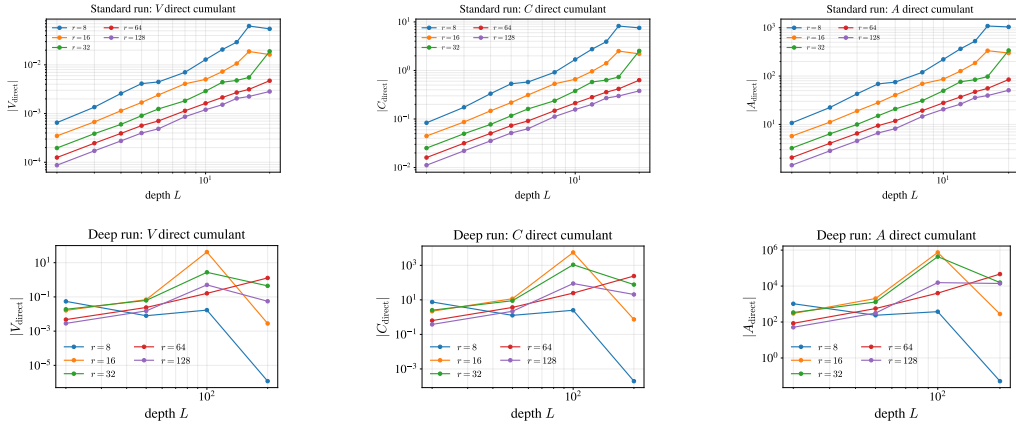


Figure 8: Log-log finite-network cumulant magnitudes. Top: the standard long run. Bottom: the $L \in \{20, 50, 100, 200\}$ deep stress test. These plots show effective depth slopes directly, before normalizing by the proposed L , L^2 , and L^3 scales.

Table 3: Exact full-rank matching for true empirical NTKs at $n = r = 128$, $m = 8$, $d = 16$. Differences remain at numerical precision up to depth $L = 200$, confirming Theorem 3.1.

Depth	max absolute NTK difference	relative difference
2	$8.88 \cdot 10^{-16}$	$8.09 \cdot 10^{-17}$
4	$3.55 \cdot 10^{-15}$	$1.31 \cdot 10^{-16}$
6	$5.33 \cdot 10^{-15}$	$4.05 \cdot 10^{-16}$
8	$3.55 \cdot 10^{-15}$	$4.20 \cdot 10^{-16}$
10	$6.22 \cdot 10^{-15}$	$1.14 \cdot 10^{-15}$
20	$1.95 \cdot 10^{-14}$	$3.77 \cdot 10^{-15}$
50	$5.68 \cdot 10^{-14}$	$6.68 \cdot 10^{-16}$
100	$2.14 \cdot 10^{-15}$	$2.35 \cdot 10^{-14}$
200	$8.08 \cdot 10^{-14}$	$1.27 \cdot 10^{-14}$

882 I MNIST Parameter-Efficiency Snapshot

883 This auxiliary experiment is included only as empirical motivation for the rank bottleneck question; it
 884 is not used in the NTK proofs. The numerical summary is reported in Table 1. MNIST images are
 885 flattened to 784 dimensions and normalized with the standard MNIST mean and variance. All models
 886 are trained with cross-entropy using Adam with learning rate 10^{-3} , batch size 128, 30 epochs, seed
 887 42, and gradient clipping at norm 1. The dense baseline is a ReLU MLP with hidden width 512 and
 888 two hidden layers. The low-rank models use the same width and a rank bottleneck; in the small-rank
 889 rows, one random feature side of the bottleneck is held fixed and the low-rank trainable maps and
 890 classifier are optimized.

891 J Why Strict Positive NMF Is Not Covered

892 This appendix explains why signed isometric low rank is not the same mathematical model as
 893 nonnegative matrix factorization. The model in (1) is signed and isometric. The right factor B
 894 is unconstrained and the frozen left factor U has orthonormal columns. Strict nonnegative matrix
 895 factorization imposes a cone constraint and changes the tangent space even at initialization. Its NTK
 896 is a projected or conic kernel, not the fixed linear-subspace kernel analyzed here. Therefore the results
 897 should not be cited as a theorem for positive NMF without a separate analysis of active constraints,
 898 boundary faces, and cone-projected gradients.

899 K Proof Details for Endpoint Constants

900 This appendix gives the elementary ReLU moment calculation that produces the endpoint constants
 901 used in the main text. For $g \sim \mathcal{N}(0, 1)$, let

$$X_0 = 2(g_+)^2, \quad Y_0 = 2\mathbf{1}_{\{g>0\}}.$$

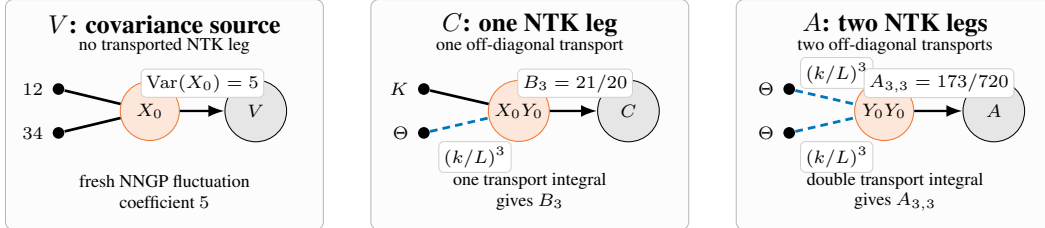


Figure 9: Endpoint diagrams for the constants in (20). Each panel separates the local ReLU atom from the external transport: V has no transported NTK leg, C has one off-diagonal NTK leg, and A has two.

902 Here X_0 is the scalar ReLU covariance atom and Y_0 is the scalar derivative atom appearing on an
 903 NTK leg. Since $\mathbb{E}[g_+^2] = 1/2$ and $\mathbb{E}[g_+^4] = 3/2$,

$$\mathbb{E}X_0 = 1, \quad \mathbb{E}X_0^2 = 6, \quad \text{Var}(X_0) = 5.$$

904 Also $\mathbb{E}Y_0 = 1$, $\mathbb{E}Y_0^2 = 2$, and therefore $\text{Var}(Y_0) = 1$. Finally,

$$\mathbb{E}X_0 Y_0 = \mathbb{E}[4(g_+)^2 \mathbf{1}_{\{g>0\}}] = 4\mathbb{E}[g_+^2] = 2, \quad \text{Cov}(X_0, Y_0) = 2 - 1 = 1.$$

905 Thus a fresh K - K source has endpoint coefficient 5.

906 It remains to explain the transported NTK legs. If a leg born at layer k is transported to depth L , the
 907 ReLU edge-of-chaos off-diagonal derivative satisfies

$$\prod_{j=k+1}^L \dot{c}(\rho_j) = \prod_{j=k+1}^L \left(1 - \frac{3}{j} + O\left(\frac{\log j}{j^2}\right) \right) = \left(\frac{k}{L}\right)^3 (1 + o(1)).$$

908 More generally, write the transport exponent as λ . The finite depth sum for one transported NTK leg
 909 is a Riemann/Beta-sum limit,

$$B_\lambda = \lim_{L \rightarrow \infty} \frac{1}{L} \sum_{k=1}^L \left(\frac{k}{L}\right)^\lambda \left(1 + 4\frac{k}{L}\right) = \frac{1}{\lambda+1} + \frac{4}{\lambda+2} = \frac{1}{\lambda+1} \left(5 - \frac{4}{\lambda+2}\right).$$

910 The two-leg calculation is the same pairing-basis projection with two transported legs. Its endpoint
 911 sum is

$$A_{\lambda,\mu} = \frac{1}{(\lambda+1)(\mu+1)} \left(5 - \frac{4}{\lambda+2} - \frac{4}{\mu+2} + \frac{4}{\lambda+\mu+3}\right).$$

912 Consequently, for the isolated off-diagonal ReLU modes used in the main text, $\lambda = \mu = 3$, and

$$V = 5, \quad C = B_3 = \frac{21}{20}, \quad A = A_{3,3} = \frac{173}{720}.$$