
Global Convergence and Better Spectral Bias in Low-Rank Neural Networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Neural networks often struggle to learn highly oscillatory functions at finite training
2 time: low-frequency components are fitted first, while high-frequency modes can
3 remain poorly recovered even when the model is expressive enough to represent
4 them. Low-rank networks are usually introduced as compressed alternatives to
5 dense models, but this view overlooks a more useful possibility: rank can act
6 as a structural control on what the network learns first. In this paper, we show
7 that this is the case. We first prove that low-rank random-feature networks in the
8 mean-field limit converge to a global minimizer of the population risk whenever
9 their limiting dynamics converge. We then show that the rank is not merely a
10 compression parameter: choosing it correctly can reduce the number of trainable
11 degrees of freedom while also improving the fit of highly oscillatory targets. The
12 key practical message is that the best rank is typically intermediate. If the rank is
13 too small, the model lacks expressivity; if it is too large, it recovers the finite-time
14 bias of the dense model. Controlled geometric and Fourier diagnostics, together
15 with high-frequency regression experiments, show that an appropriate low rank can
16 lower test loss, improve high-frequency recovery, and that the optimal rank shifts
17 with the target spectrum and training objective.

18 1 Introduction

19 Low-rank neural networks replace dense matrices by factored or bottlenecked operators. The standard
20 motivation is computational: fewer trainable parameters and cheaper memory traffic. This paper
21 argues that the more important question is mathematical. Can rank be used as a principled control on
22 training dynamics rather than only as a compression parameter?

23 This question is important for scientific machine learning. Many tasks in AI for science require
24 learning functions with oscillations, sharp transitions, multiscale structure, or high-frequency Fourier
25 content: wave fields, PDE solution operators, molecular potentials, turbulent signals, and inverse
26 problems all contain information at small spatial or temporal scales. Neural networks can represent
27 such functions, but finite-time training often learns smooth, low-frequency components first. This
28 frequency bias can be modified by the data distribution or by replacing the loss with a Sobolev
29 norm [23]; our question is whether the architecture itself, through rank, can provide another control
30 parameter. This makes the problem difficult: expressivity alone is not enough, and a model that is
31 globally trainable in principle can still fail to recover the high-frequency part under a realistic training
32 budget.

33 We answer the first question positively for a low-rank random-feature architecture. The model
34 freezes a rich random-feature map and trains low-rank channel weights through mean-field gradient
35 flow. Since the frozen features keep a dense span throughout training, the usual obstruction in deep
36 mean-field convergence proofs is removed: at a limit point, zero gradient implies the conditional

37 first-order condition for the loss, and the loss-identifiability assumption forces global optimality. Low
38 rank enters through the mixing matrix and channel index, changing constants but not the logic of the
39 proof.

40 The second question is more subtle. Global convergence is an asymptotic statement: it says what
41 happens if the population dynamics settle. Spectral bias is a finite-time phenomenon: it asks which
42 parts of the target are reached first under a limited optimization budget. These two views can disagree
43 on oscillatory targets, because a model may be globally trainable and still spend most of training
44 on the low-frequency part of the signal. We show that rank changes this finite-budget behavior by
45 changing the geometry of ReLU networks, hence which Fourier modes are recovered early. The
46 practical outcome is a rank-selection rule: choose the smallest rank that avoids the approximation
47 bottleneck while improving high-frequency recovery.

48 **Contributions.** The paper makes three contributions:

- 49 • We give a globally convergent anchor model for low-rank networks. Freezing a rich random-
50 feature map keeps the feature span from collapsing, while the trainable low-rank channels
51 still allow a mean-field feature-learning interpretation.
- 52 • We explain why rank can change high-frequency learning. In ReLU networks, lowering rank
53 reduces the number of independent path constraints, which changes the piecewise-affine
54 geometry that controls Fourier energy.
- 55 • We provide controlled diagnostics for choosing rank. The experiments show that the useful
56 rank is usually intermediate and moves with the target spectrum or training objective, so
57 rank should be selected from spectral recovery rather than parameter count alone.

58 2 Related Work

59 This section positions the paper relative to mean-field convergence, low-rank adaptation, bottleneck-
60 rank theory, and spectral bias.

61 **Mean-field convergence.** Mean-field theory analyzes neural network training through the deter-
62 ministic evolution of parameter distributions in the large-width limit. Foundational work established
63 global convergence mechanisms for two-layer models and Wasserstein gradient flows under suitable
64 convexity, homogeneity, or support assumptions [17, 5]. Later multilayer mean-field analyses clarified
65 how deep feature maps evolve and when stationary points can be related to global optima [20, 19].
66 Our result follows this line but focuses on low-rank random-feature networks: the existence and
67 uniqueness estimates use the same mean-field template, while the convergence proof must be adapted
68 to channel-wise low-rank couplings.

69 **Low-rank and random features.** Low-rank neural networks are often motivated by parameter
70 efficiency, with LoRA-style decompositions now standard in large-scale adaptation [15]. Recent
71 theory has begun to explain why LoRA can be trainable: Jang, Lee, and Ryu [16] show that, in the
72 NTK fine-tuning regime, sufficiently large LoRA rank removes spurious local minima and yields low-
73 rank solutions with good generalization. Most theory for low-rank networks focuses on expressivity,
74 approximation, kernel behavior, or LoRA fine-tuning rather than global mean-field training dynamics.
75 Random features provide a tractable bridge between finite networks and kernel methods: fixing the
76 feature map turns the representation into a stable basis, while training the channel weights retains
77 a mean-field optimization problem. In our RF-LR setting, frozen random features prevent support
78 collapse and preserve the dense-span property needed for the stationarity-to-optimality argument; the
79 low-rank factors enter through bounded channel mixing rather than through a full dense contraction
80 chain.

81 **Bottleneck rank and implicit low-rank bias.** Recent work gives a functional notion of optimal rank
82 that is closer to our setting than the algebraic rank of a single matrix. Jacot [12] studies the implicit
83 bias of large-depth homogeneous networks and shows that their representation cost converges toward
84 a nonlinear rank notion, sandwiched between Jacobian rank and bottleneck rank. The bottleneck
85 rank is the smallest inner dimension k for which a function can be factored as $f = h \circ g$ through
86 \mathbb{R}^k . Follow-up work [13] proves that learned deep representations exhibit a bottleneck structure

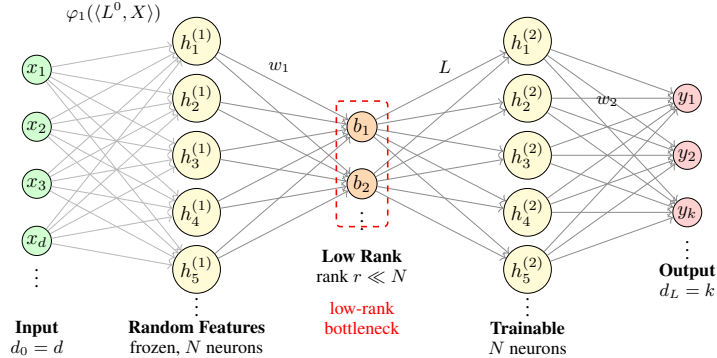


Figure 1: Architecture of a three-layer low-rank random-feature network with frozen random features, a rank- r bottleneck, and trainable output weights.

87 and introduces finite-depth corrections that balance low inner dimension against regularity. This
 88 suggests a theoretical target rank: the smallest rank that captures the intrinsic bottleneck dimension
 89 of the target without forcing an irregular or rank-underestimating interpolant. Complementarily,
 90 Bantzis, Simon, and Jacot [3] show that the first saddle escape of deep ReLU networks has a low-rank
 91 bias in deeper layers and propose saddle-to-saddle dynamics with increasing bottleneck rank. Our
 92 rank-selection principle is consistent with this picture: training should start from low effective rank,
 93 but useful learning requires increasing rank until approximation and spectral recovery saturate.

94 **Spectral bias.** The spectral-bias literature identifies why neural networks often fit low-frequency
 95 components before high-frequency components. Rahaman et al. [21] connect this learning order to
 96 Fourier properties of piecewise linear networks, while Yu, Yang, and Townsend [23] show that the
 97 bias can be tuned through nonuniform data and Sobolev-type losses. This paper keeps that Fourier
 98 viewpoint but changes the intervention: instead of reweighting data or losses, we study rank as an
 99 architectural parameter that changes the geometry seen by finite-time training.

100 3 Freezing Half the Weights Allows Global Convergence

101 This section defines the RF-LR mean-field model and states the conditional global convergence
 102 theorem.

103 Let $X \in \mathbb{R}^d$ and Y be drawn from a population distribution. The finite-width low-rank random-
 104 feature network used throughout the paper is the one from the original experiments:

$$h^{(0)}(x) = x, \quad h^{(\ell)}(x) = \frac{1}{n_\ell} \sum_{j=1}^{n_\ell} w_j^{(\ell)} \varphi_\ell \left(\langle L_j^{(\ell)}, h^{(\ell-1)}(x) \rangle + b_j^{(\ell)} \right), \quad (1)$$

105 where the feature directions $L_j^{(\ell)}$ and biases $b_j^{(\ell)}$ are frozen, while the channel weights $w_j^{(\ell)}$ are trained.
 106 This is a low-rank factorization of a dense layer in which the left factor is fixed: if $W_\ell = U_\ell V_\ell^\top$ with
 107 $U_\ell \in \mathbb{R}^{n_\ell \times r_\ell}$ bounded and full column rank, then U_ℓ plays the role of a frozen random-feature or
 108 mixing basis and the trainable factor V_ℓ is absorbed into the channel weights. Thus one may avoid
 109 training U_ℓ entirely; drawing it as a bounded full-rank tall-skinny matrix is enough for the theory,
 110 while choosing $U_\ell \in \text{St}(n_\ell, r_\ell)$ is an optional numerical normalization.

111 3.1 Mean-Field Formulation

112 The mean-field viewpoint is the relevant one for global convergence. Pure kernel or neural-tangent
 113 analyses give convexity by freezing the representation, but they largely remove feature learning. Finite-
 114 width nonconvex analyses allow feature learning but rarely provide global convergence guarantees
 115 for deep low-rank networks. Mean-field theory is currently the most developed framework that can
 116 keep a feature-learning interpretation while still proving that limiting gradient-flow trajectories reach
 117 global minimizers under identifiability assumptions. We use a frozen-feature low-rank specialization

118 because it is the cleanest setting in which the dense-span argument is preserved; it should be viewed
 119 as the globally convergent anchor case for more general low-rank feature-learning models.

120 In the mean-field limit, frozen feature maps are written $L^0(C_\ell)$ and trainable channel weights are
 121 written $w_\ell(t, C_\ell)$. The channel mixer $L \in \mathbb{R}^{N_2 \times r}$ is a bounded full-column-rank tall-skinny factor,
 122 not necessarily an orthogonal matrix. We write the second-layer operator as $W_2 = LA_2^\top$, with L
 123 fixed and A_2 trainable, then absorb the small right factor A_2^\top into the rank-channel functions below.
 124 One may additionally choose $L \in \text{St}(N_2, r)$, or periodically re-orthogonalize it, as a numerical
 125 design choice; this is natural for matrix-structured optimizers such as Muon, which orthogonalize
 126 matrix updates by Newton–Schulz or polar iterations [14, 4]. In the two-hidden-layer case, the first
 127 partial functions are

$$f_k(t, x) = \mathbb{E}_{C_1} [w_1(t, C_1, k) \varphi_1((L^0(C_1), x))], \quad k = 1, \dots, r, \quad (2)$$

128 and the second-layer preactivation is

$$H_2(t, c_2; x) = \sum_{k=1}^r L_{c_2, k} f_k(t, x). \quad (3)$$

129 Thus H_2 is the reconstruction of the second-layer signal from a bounded rank- r channel basis. The
 130 output is obtained by averaging the trainable top weights against $\varphi_2(H_2)$. For deeper networks,
 131 the same alternating pattern is used: trainable channel weights are averaged against frozen random
 132 features and mixed through bounded low-rank channel matrices.

133 The population objective is

$$\mathcal{L}(W) = \mathbb{E}_{(X, Y)} [\mathcal{L}(Y, \hat{y}(X; W))]. \quad (4)$$

134 The mean-field gradient flow evolves the trainable weights by

$$\partial_t w_\ell(t, \cdot) = -\xi_\ell(t) \nabla_{w_\ell} \mathcal{L}(W(t)), \quad (5)$$

135 where ξ_ℓ are bounded learning-rate schedules. The exact formulas are standard backpropagation in
 136 the mean-field variables; the only low-rank modification is the channel mixing through $L_{c_\ell, k}$.

137 Under the technical regularity, diversity, loss-identifiability, and convergence assumptions stated in
 138 Appendix A, freezing the feature directions gives a well-defined mean-field system and allows the
 139 usual stationarity argument to conclude global optimality.

140 **Theorem 3.1** (Well-posedness and conditional global convergence of RF-LR training). *Under*
 141 *Assumptions A.1–A.6, the low-rank mean-field gradient-flow system has a unique solution on every*
 142 *finite interval $[0, T]$ and therefore on $[0, \infty)$. Moreover, if the resulting low-rank mean-field dynamics*
 143 *converges to W^* in the modified \mathcal{W}_4 coupling topology defined in Appendix A through Assumption A.6,*
 144 *then W^* is a global minimizer of the population loss:*

$$\mathcal{L}(W^*) = \inf_W \mathcal{L}(W). \quad (6)$$

145 *For every depth $L \geq 2$, this holds with standard independent initialization of the trainable weights.*

146 The convergence statement is conditional: as in the corresponding multilayer mean-field results, we
 147 do not prove that every trajectory converges. The theorem says that convergence cannot occur to a
 148 bad stationary point under the stated dense-span, non-degeneracy, and loss-identifiability assumptions.
 149 The well-posedness part only justifies that the limiting dynamics are genuine ODEs rather than formal
 150 flows. This is the natural Picard argument of Nguyen and Pham [19]: nothing structural changes in
 151 the proof, except that dense-layer operator norms are replaced by bounded low-rank mixing constants.
 152 The complete well-posedness proof is deferred to the supplementary material, Appendix B.

153 *Proof idea.* The existence and uniqueness part follows the standard mean-field fixed-point argument,
 154 whose full details are given in Appendix B. We only summarize the convergence-to-global-minimum
 155 step here; Appendix C gives the full notation and the five formal steps. We assume that gradient
 156 descent, or its mean-field gradient-flow limit, has converged somewhere: $W(t) \rightarrow \bar{W}$ in the modified
 157 \mathcal{W}_4 coupling topology defined in Appendix A. Stationarity and the persistent dense span of the frozen
 158 first-layer features then imply the integrated identity

$$\mathbb{E} \left[d_L(Z; \bar{W}) B_k^{(2)}(X; \bar{W}) \mid X = x \right] = 0 \quad \text{for } \mathbb{P}_X\text{-a.e. } x, \quad k = 1, \dots, r, \quad (7)$$

159 where $d_L(Z; \bar{W}) = \partial_2 \mathcal{L}(Y, \hat{y}(X; \bar{W}))$ and $B_k^{(2)}$ is the channel-wise backpropagated factor defined
 160 in Appendix C.

161 *Step 3: remove the backpropagated factor.* The non-degeneracy assumptions exclude a dead-channel
 162 limit: for \mathbb{P}_X -a.e. x , at least one channel has a nonzero backpropagated factor $B_k^{(2)}(x; \bar{W})$. The
 163 identity we use keeps this factor inside the conditional expectation:

$$\mathbb{E}\left[d_L(Z; \bar{W})B_k^{(2)}(X; \bar{W}) \mid X = x\right] = 0 \quad \text{for } \mathbb{P}_X\text{-a.e. } x, \quad k = 1, \dots, r.$$

164 Since $B_k^{(2)}(X; \bar{W})$ is X -measurable, the non-degeneracy assumption then allows us to choose at
 165 least one nonzero backpropagated factor at almost every x , so

$$\mathbb{E}\left[\partial_2 \mathcal{L}(Y, \hat{y}(X; \bar{W})) \mid X = x\right] = 0 \quad \text{for } \mathbb{P}_X\text{-a.e. } x. \quad (8)$$

166 *Step 5: the trajectory loss converges to the limit loss.* Since $W(t) \rightarrow \bar{W}$ in the modified \mathcal{W}_4 coupling
 167 topology defined in Appendix A, there are couplings π_t between $W(t)$ and \bar{W} whose weighted
 168 channel gaps vanish. In the two-layer channel notation, the key convergence statement is exactly

$$\mathbb{E}_{\pi_t} \left[(1 + |\bar{w}_2|) |\bar{w}_2| \sum_{k=1}^r |\bar{w}_{1,k}| |w_{1,k} - \bar{w}_{1,k}| \right] \rightarrow 0, \quad (9)$$

169 with analogous quantities for all deeper layers. The forward stability estimate uses the ReLU channel
 170 expansion of a low-rank preactivation,

$$\left| \text{ReLU} \left(\sum_{k=1}^r x_k \right) - \text{ReLU} \left(\sum_{k=1}^r y_k \right) \right| \leq \sum_{k=1}^r |x_k - y_k|. \quad (10)$$

171 Together with the low-rank form $H_2 = \sum_k L_{c_2, k} f_k$, (10) bounds every forward and backward
 172 difference by the same channel integrals. Hence

$$\mathbb{E} \left[|\hat{y}(X; W(t)) - \hat{y}(X; \bar{W})| \right] \rightarrow 0. \quad (11)$$

173 Lipschitzness of the loss in its second argument on the a priori bounded trajectory yields

$$|\mathcal{L}(W(t)) - \mathcal{L}(\bar{W})| \leq K \mathbb{E} \left[|\hat{y}(X; W(t)) - \hat{y}(X; \bar{W})| \right] \rightarrow 0. \quad (12)$$

174 *Step 4: identify the limit as globally optimal.* By the loss-identifiability condition in Assumption A.3,
 175 namely (26) in Appendix A,

$$\mathbb{E}[\partial_2 \mathcal{L}(Y, u) \mid X = x] = 0 \quad \implies \quad \mathbb{E}[\mathcal{L}(Y, u) \mid X = x] = 0. \quad (13)$$

176 Applying this with $u = \hat{y}(x; \bar{W})$ and using non-negativity of \mathcal{L} gives

$$\mathcal{L}(\bar{W}) = 0 = \inf_W \mathcal{L}(W). \quad (14)$$

177 Together with Step 5, $\mathcal{L}(W(t)) \rightarrow \mathcal{L}(\bar{W})$. □

178 4 Low-Rank Spectral Bias

179 This section explains how low-rank path constraints change CPWL geometry and the Fourier coeffi-
 180 cients that control finite-resolution spectral bias.

181 **Connection to spectral bias.** Theorem 3.1 says that a convergent trajectory reaches a global
 182 minimizer, but it does not describe the order in which Fourier modes are recovered at finite time.
 183 This section explains the architectural mechanism behind that order: low rank constrains active
 184 path coefficients and switching geometry, so it changes the CPWL face moments that control finite-
 185 resolution Fourier transfer. Thus the role of Rahaman et al. [21] here is not only as prior related work,
 186 but as the Fourier template that our rank-dependent face-moment calculation refines.

187 In a ReLU network, on every activation region the network is affine, and the affine coefficient is a
 188 sum over active paths. We use the same path-product viewpoint as in path-based analyses of loss

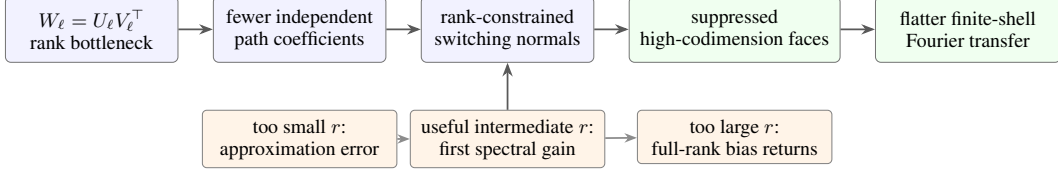


Figure 2: Mechanism summary linking low-rank factorization, path coefficients, switching normals, CPWL faces, and finite-shell Fourier transfer.

189 surfaces and finite-width neural kernels [6, 10]. For a scalar-output depth- L network, this can be
190 written schematically as

$$f(x) = \sum_{p=(i_0, \dots, i_L)} \left(\prod_{\ell=1}^L W_{i_\ell i_{\ell-1}}^{(\ell)} \right) \left(\prod_{\ell=1}^{L-1} \mathbf{1}_{\{h_{i_\ell}^{(\ell)}(x) > 0\}} \right) x_{i_0}, \quad (15)$$

191 where p ranges over input-to-output paths. If $W_\ell = U_\ell V_\ell^\top$ has rank r_ℓ , then many symbolic paths
192 share the same latent channels, so their coefficients cannot vary independently. These path constraints
193 also constrain ReLU switching normals $\nabla_x h_j^{(\ell)}(x)$: the normals lie in a smaller span than in a dense
194 layer. Geometrically, this suppresses low-dimensional, high-codimension intersections; spectrally,
195 the Fourier shell law below predicts more visible facet-dominated high-frequency mass.

196 **Mechanism summary.** Figure 2 summarizes the mechanism. A low-rank factorization reduces
197 independent path coefficients in (15); this constrains switching normals, suppresses high-codimension
198 CPWL faces, and moves shell-averaged Fourier energy toward the facet-dominated regime of Ra-
199 haman et al. [21]. The useful rank is intermediate: too small gives approximation error, while too
200 large recovers the full-rank low-frequency-first bias. A frequency shell means modes with comparable
201 radius $\|\xi\| \simeq \rho$.

202 This connects to bottleneck rank: the population lower bound is the smallest inner dimension through
203 which the target can factor [12, 13], while the finite-budget choice also balances optimization and
204 spectral recovery. In experiments, we estimate this by the first rank where approximation is controlled
205 and Fourier recovery no longer improves.

206 This algebraic constraint becomes geometric. Inside a cell of the previous layers, a new switching
207 normal has the form

$$n_{\ell i \varepsilon} = B_{\ell-1, \varepsilon}^\top V_\ell u_{\ell i} \in \text{Im}(B_{\ell-1, \varepsilon}^\top V_\ell), \quad \dim \text{Im}(B_{\ell-1, \varepsilon}^\top V_\ell) \leq \min(d, r_1, \dots, r_\ell). \quad (16)$$

208 Thus rank lowers the dimension of the normal subspace and suppresses high-codimension faces first.

209 **Face Fourier expansion.** Let $g = \chi f$, where f is CPWL on a finite polyhedral complex in \mathbb{R}^d and
210 $\chi \in C_c^\infty(\mathbb{R}^d)$ is an observation window. Let $\mathcal{F}_q(f)$ be the codimension- q faces of the complex. For
211 every large frequency $\xi = \rho\theta$, $\theta \in \mathbb{S}^{d-1}$, the Fourier transform has the expansion

$$\widehat{g}(\rho\theta) = \sum_{q=1}^d \rho^{-(q+1)} \sum_{F \in \mathcal{F}_q(f)} A_F(\rho, \theta) + O(\rho^{-(d+2)}), \quad (17)$$

212 where A_F is a bounded oscillatory face coefficient controlled by the jump Δ_F of the appropriate
213 normal derivative across the local fan around F . In particular,

$$|A_F(\rho, \theta)| \leq C_\chi \|\Delta_F\| \text{vol}_{d-q}(F) \omega_F(\theta). \quad (18)$$

214 This is the anisotropic phenomenon isolated by Rahaman et al. [21]: in almost all directions a ReLU
215 network has fast $k^{-(d+1)}$ amplitude decay, while in special directions orthogonal to facets of the
216 linear regions the decay can be as slow as k^{-2} . Thus low codimension faces are precisely the
217 geometric structures that keep high-frequency Fourier mass visible at finite resolution.

218 **Proposition 4.1** (Fourier shell law for CPWL networks). *Let $g = \chi f$ be a windowed continuous*
219 *piecewise affine network in dimension d . Define the codimension- q face moment*

$$M_q(f) = \sum_{F \in \mathcal{F}_q(f)} \|\Delta_F\|^2 \text{vol}_{d-q}(F)^2 \mathbb{E}_{\theta \in \mathbb{S}^{d-1}} \omega_F(\theta)^2. \quad (19)$$

220 Then the shell-averaged Fourier energy satisfies

$$S_f(\rho) = \mathbb{E}_{\|\xi\| \simeq \rho} |\widehat{g}(\xi)|^2 \lesssim \sum_{q=1}^d M_q(f) \rho^{-2(q+1)}. \quad (20)$$

221 Low rank suppresses M_q for larger q first, moving finite-resolution spectra toward the facet-dominated
 222 ρ^{-4} energy regime. This should be read as a coefficient statement rather than a claim that low rank
 223 creates a new universal exponent: the useful comparison with full rank is the reduction of M_q/M_1
 224 and of the transition coefficient $(M_q/M_1)\rho^{-2q+2}$ for $q \geq 2$.

225 *Proof idea.* Write f on its polyhedral complex as an affine function on each cell and set $g = \chi f$ with
 226 a smooth compactly supported window. The Fourier transform is then a sum of ordinary oscillatory
 227 integrals over cells. Integrating by parts in the direction $\theta = \xi/\|\xi\|$ produces boundary terms on cell
 228 facets. The order- ρ^{-1} traces cancel across interior facets because the CPWL network is continuous
 229 and the same window multiplies both sides.

230 The first non-canceling term is therefore the jump of the directional derivative across a facet, hence it
 231 carries the factor ρ^{-2} . Repeating the same integration-by-parts argument on the induced face complex
 232 localizes the next non-smooth terms on codimension- q faces, with amplitude order $\rho^{-(q+1)}$. This
 233 gives the face expansion (17), with coefficients controlled by the corresponding normal-derivative
 234 jumps and face volumes. This is the same anisotropic mechanism emphasized by Rahaman et
 235 al. [21]: generic directions decay quickly, while directions close to facet normals can retain slow
 236 $|\xi|^{-2}$ amplitude decay.

237 Squaring the expansion and averaging over a frequency shell gives diagonal face energies of order
 238 $M_q(f)\rho^{-2(q+1)}$. Oscillatory cross terms are either lower order by angular non-stationary phase or
 239 absorbed into the same moments by Cauchy–Schwarz. Thus the exponent is determined by codimen-
 240 sion, while rank changes the coefficients $M_q(f)$ by suppressing high-codimension intersections. A
 241 detailed integration-by-parts proof is given in Appendix D. \square

242 5 Experiments

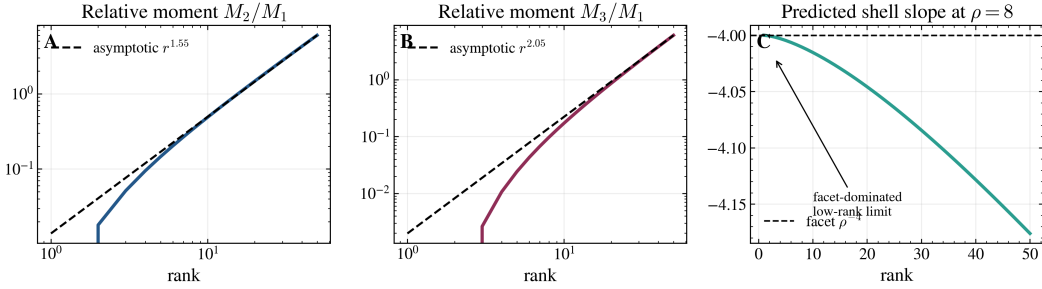


Figure 3: Theory proxy for the Fourier shell law. Panels A–B show relative face-moment proxies M_2/M_1 and M_3/M_1 , i.e. how much higher-codimension geometry is present relative to facets. Panel C shows the induced finite-shell slope, separating the universal codimension exponents from the rank-dependent coefficients.

243 This section presents controlled diagnostics that test the proposed geometric and Fourier mechanism,
 244 together with supporting high-frequency regression evidence.

245 We keep the experimental evidence deliberately short and defer the full setup of every generated plot
 246 to Appendix E. The main figures are controlled CPWL and Fourier/kernel-limit diagnostics: they test
 247 the mechanism predicted by Proposition 4.1, namely that rank changes the face-moment coefficients
 248 $M_q(f)$ and the finite-time recovery of Fourier modes. Separate finite-network MMNN sweeps on
 249 high-frequency targets are used only as supporting evidence that low-rank models can fit oscillatory
 250 functions and that optimizer choice and feature freezing matter in practice.

251 The experiments should therefore be read as mechanism checks rather than as a benchmark claim
 252 that one optimizer or one rank universally dominates. The CPWL plots isolate geometry, the

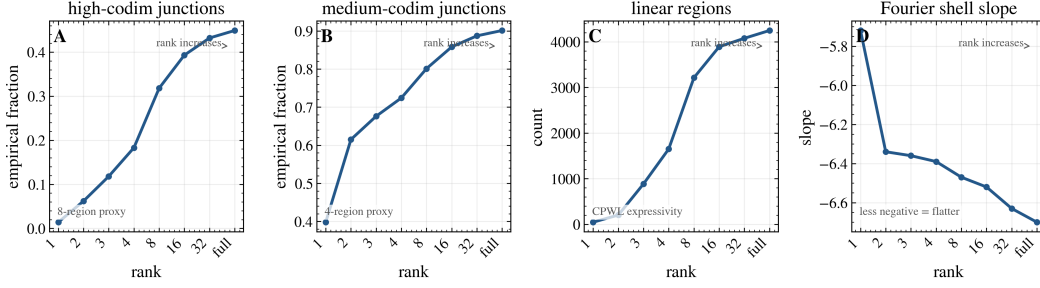


Figure 4: Three-dimensional CPWL geometry proxy. Panels A–B show the fraction of sampled neighborhoods where many affine regions meet, using thresholds of 8 and 4 regions as empirical high- and medium-codimension junction proxy. Panel C shows the number of linear regions, and Panel D shows the estimated shell slope.

253 Fourier/kernel plots isolate spectral transfer at fixed budget, and the rank-shift control tests whether
 254 the observed optimum is task-dependent. This separation is intentional: the theorem is an asymptotic
 255 convergence statement, whereas the spectral-bias claims concern finite-resolution coefficients and
 256 finite-time recovery.

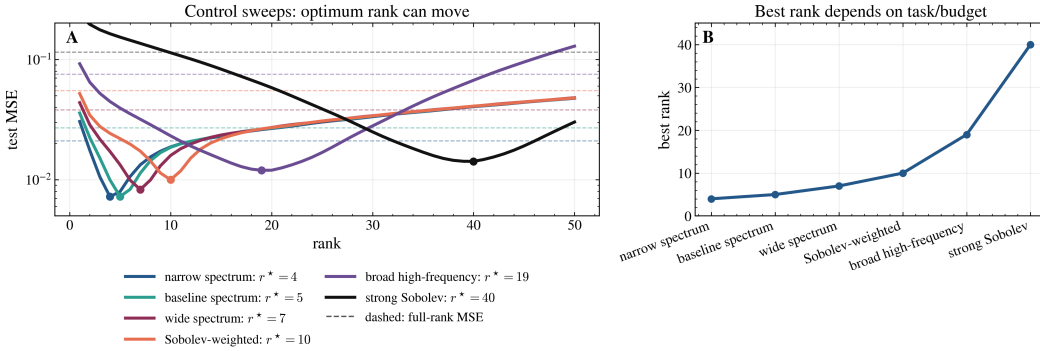


Figure 5: Rank-shift control. The left panel shows finite-budget MSE versus rank for several target/weighting regimes; the right panel shows the best rank in each regime. Dashed lines are full-rank baselines. The plot emphasizes that the useful rank moves with the spectrum and objective rather than staying fixed.

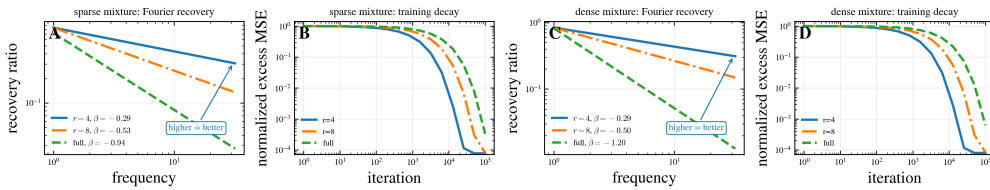


Figure 6: Mode-wise recovery and training curves. Panels A and C show Fourier recovery ratios $R_r(k, t)$ across frequency modes: higher and flatter curves mean that high frequencies are recovered more evenly relative to low frequencies. Panels B and D show normalized excess MSE during training, so the useful ranks are those that combine low error with flatter spectral recovery.

257 The one-dimensional recovery plots require a small caveat. A 1D ReLU network is a linear spline: if
 258 Δ_j denotes the derivative jump at knot t_j , then

$$\hat{f}(k) = -\frac{1}{(ik)^2} \sum_j \Delta_j e^{-ikt_j} + O(|k|^{-3}), \quad |\hat{f}(k)|^2 \sim |k|^{-4} A(k). \quad (21)$$

259 Thus the raw asymptotic tail stays near the classical k^{-4} law, and the finite-budget diagnostic should
 260 be mode-wise recovery,

$$R_r(k, t) = \frac{|\widehat{f}_{r,t}(k)|}{|\widehat{f}^*(k)|}, \quad \beta_r(t) = \frac{d \log R_r(k, t)}{d \log k}. \quad (22)$$

261 A flatter fitted slope $\beta_r(t)$ means that high modes are recovered more evenly relative to low modes
 262 on the plotted frequency window. Figure 6 reports this recovery ratio together with training MSE.

263 The geometry plots should be read through the face-moment shell law

$$S_f(\rho) \lesssim M_1(f)\rho^{-4} + M_2(f)\rho^{-6} + M_3(f)\rho^{-8} + \dots \quad (23)$$

264 Figure 3 tests the coefficient prediction directly: higher-codimension moments grow with rank,
 265 so the finite-shell slope moves away from the facet-dominated regime. Figure 4 checks the same
 266 mechanism geometrically: high-codimension junctions are empirical proxies for the higher moments
 267 M_2, M_3, \dots , and they become more frequent as rank creates more independent switching geometry.
 268 Appendix E gives the full interpretation.

269 The spectral-bias experiment isolates the rank effect in a width- 2^{15} Fourier/kernel surrogate. Figure 5
 270 is the key control: the best finite-budget rank is intermediate, but it is not fixed at one value. It moves
 271 from small ranks on narrow or baseline spectra to larger ranks under wider spectra and Sobolev-
 272 weighted objectives used to tune frequency bias [23]. Across these regimes, the best low-rank point
 273 also beats the corresponding full-rank dashed baseline. Figure 6 shows why MSE alone is not enough:
 274 the useful rank also keeps high-frequency recovery flatter than the full endpoint.

275 Together, Figures 3–6 support the same conclusion: rank is not only a compression parameter. It
 276 changes the finite-budget spectral transfer, and the best rank is the first one that controls approximation
 277 while preserving flatter high-frequency recovery.

278 6 Conclusion

279 Low rank is not only a compression device. In the random-feature mean-field regime, it preserves
 280 global convergence because frozen features maintain dense span and low-rank channel mixing only
 281 changes constants. In ReLU and CPWL regimes, it also changes spectral bias through path-space
 282 constraints and Fourier face moments. This matters for AI for science because oscillatory targets are
 283 common in PDEs, waves, inverse problems, and molecular or multiscale modeling. Prior work shows
 284 that frequency bias can be tuned by data or Sobolev losses [23]; our contribution is to show that
 285 architecture provides another lever. The message is conservative but actionable: use the frozen-feature
 286 low-rank model as a globally convergent anchor, then choose rank by finite-budget spectral recovery
 287 rather than by parameter count alone.

288 The mechanism is the following. A globally convergent dynamic can still learn low frequencies first
 289 and fail at finite time on highly oscillatory targets. Low rank changes this finite-time behavior because,
 290 in ReLU networks, it constrains active path coefficient tensors; these constraints restrict switching
 291 normals, suppress high-codimension faces in the continuous piecewise affine complex, and change
 292 Fourier shell behavior. Here CPWL means continuous piecewise affine: the input space is split into
 293 polyhedral regions, and the network is affine on each region. In one dimension, the raw spline tail
 294 remains close to the classical k^{-4} law, so the right diagnostic is not the asymptotic exponent alone but
 295 the learned Fourier recovery curve. This leads to the practical rule used throughout the experiments:
 296 choose the smallest rank that avoids the approximation bottleneck while flattening the finite-time
 297 spectral transfer.

298 Several directions remain open. First, the convergence result should be extended from frozen
 299 random features to trainable low-rank factors, where the dense-span argument can fail if the features
 300 collapse. Second, the CPWL shell law suggests a measurable rank-dependent coefficient theory,
 301 but sharper finite-width estimates are needed to predict the best rank without a sweep. Third, the
 302 experiments indicate that optimizer geometry matters, especially for Muon-type orthogonalized
 303 updates; understanding how such updates interact with low-rank spectral transfer is a natural next
 304 step.

References

- [1] F. Bach. Breaking the curse of dimensionality with convex neural networks. *Journal of Machine Learning Research*, 18(19):1–53, 2017. arXiv:1412.8690, <https://arxiv.org/abs/1412.8690>.
- [2] A. R. Barron. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information Theory*, 39(3):930–945, 1993. doi:10.1109/18.256500.
- [3] I. Bantzis, J. B. Simon, and A. Jacot. Saddle-to-saddle dynamics in deep ReLU networks: low-rank bias in the first saddle escape. arXiv preprint arXiv:2505.21722, 2026. <https://arxiv.org/abs/2505.21722>.
- [4] V. Boreiko, Z. Bu, and S. Zha. Towards understanding orthogonalization in Muon. In *HiLD Workshop at ICML*, 2025. OpenReview: <https://openreview.net/forum?id=ppmyFtr9EW>.
- [5] L. Chizat and F. Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. In *Advances in Neural Information Processing Systems 31*, pages 3036–3046, 2018. arXiv:1805.09545, <https://arxiv.org/abs/1805.09545>.
- [6] A. Choromanska, M. Henaff, M. Mathieu, G. Ben Arous, and Y. LeCun. The loss surfaces of multilayer networks. In *Proceedings of the 18th International Conference on Artificial Intelligence and Statistics*, PMLR 38:192–204, 2015. arXiv:1412.0233, <https://arxiv.org/abs/1412.0233>.
- [7] W. Czarnecki, S. Osindero, M. Jaderberg, G. Swirszcz, and R. Pascanu. Sobolev training for neural networks. In *Advances in Neural Information Processing Systems 30*, pages 4278–4287, 2017. arXiv:1706.04859, <https://arxiv.org/abs/1706.04859>.
- [8] R. Diaz, Q.-N. Le, and S. Robins. Fourier transforms of polytopes, solid angle sums, and discrete volume. arXiv preprint arXiv:1602.08593, 2018. <https://arxiv.org/abs/1602.08593>.
- [9] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics*, PMLR 9:249–256, 2010. <https://proceedings.mlr.press/v9/glorot10a.html>.
- [10] B. Hanin and M. Nica. Finite depth and width corrections to the neural tangent kernel. In *International Conference on Learning Representations*, 2020. arXiv:1909.05989, <https://arxiv.org/abs/1909.05989>.
- [11] A. Jacot, F. Gabriel, and C. Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems 31*, 2018. arXiv:1806.07572, <https://arxiv.org/abs/1806.07572>.
- [12] A. Jacot. Implicit bias of large depth networks: a notion of rank for nonlinear functions. arXiv preprint arXiv:2209.15055, 2023. <https://arxiv.org/abs/2209.15055>.
- [13] A. Jacot. Bottleneck structure in learned features: low-dimension vs regularity tradeoff. arXiv preprint arXiv:2305.19008, 2024. <https://arxiv.org/abs/2305.19008>.
- [14] K. Jordan, Y. Jin, V. Boza, J. You, F. Cesista, L. Newhouse, and J. Bernstein. Muon: An optimizer for hidden layers in neural networks. Blog post, 2024. <https://kellerjordan.github.io/posts/muon/>.
- [15] E. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. arXiv:2106.09685, <https://arxiv.org/abs/2106.09685>.
- [16] U. Jang, J. D. Lee, and E. K. Ryu. LoRA training in the NTK regime has no spurious local minima. In *Proceedings of the 41st International Conference on Machine Learning*, PMLR 235:21306–21328, 2024. arXiv:2402.11867, <https://arxiv.org/abs/2402.11867>.

- 351 [17] S. Mei, A. Montanari, and P.-M. Nguyen. A mean field view of the landscape of two-layer
352 neural networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671,
353 2018. doi:10.1073/pnas.1806579115; arXiv:1804.06561.
- 354 [18] G. Montufar, R. Pascanu, K. Cho, and Y. Bengio. On the number of linear regions of deep
355 neural networks. In *Advances in Neural Information Processing Systems 27*, pages 2924–2932,
356 2014. arXiv:1402.1869, <https://arxiv.org/abs/1402.1869>.
- 357 [19] P.-M. Nguyen and H. T. Pham. A rigorous framework for the mean-field limit of multilayer
358 neural networks. *Mathematical Statistics and Learning*, 6(3):201–357, 2023. arXiv:2001.11443,
359 <https://arxiv.org/abs/2001.11443>.
- 360 [20] H. T. Pham and P.-M. Nguyen. Global convergence of three-layer neural networks in the mean
361 field regime. In *International Conference on Learning Representations*, 2021. arXiv:2105.05228,
362 <https://arxiv.org/abs/2105.05228>.
- 363 [21] N. Rahaman, A. Baratin, D. Arpit, F. Draxler, M. Lin, F. Hamprecht, Y. Bengio, and A.
364 Courville. On the spectral bias of neural networks. In *Proceedings of the 36th International
365 Conference on Machine Learning*, PMLR 97:5301–5310, 2019. arXiv:1806.08734, <https://arxiv.org/abs/1806.08734>.
- 367 [22] T. Zaslavsky. Facing up to arrangements: face-count formulas for partitions of space
368 by hyperplanes. *Memoirs of the American Mathematical Society*, 1(154):1–102, 1975.
369 doi:10.1090/memo/0154.
- 370 [23] A. Yu, Y. Yang, and A. Townsend. Tuning frequency bias in neural network training with
371 nonuniform data. arXiv preprint arXiv:2205.14300, 2022. [https://arxiv.org/abs/2205.](https://arxiv.org/abs/2205.14300)
372 14300.

373 NeurIPS Paper Checklist

374 1. Claims

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

377 Answer: [Yes].

378 Justification: The abstract and introduction state the conditional nature of the convergence
379 result, the architectural scope, and the finite-budget spectral-bias claim. The experiments
380 are described as controlled diagnostics rather than universal benchmarks.

381 Guidelines:

- 382 • The answer [N/A] means that the abstract and introduction do not include the claims
383 made in the paper.
- 384 • The abstract and/or introduction should clearly state the claims made, including the
385 contributions made in the paper and important assumptions and limitations. A [No] or
386 [N/A] answer to this question will not be perceived well by the reviewers.
- 387 • The claims made should match theoretical and experimental results, and reflect how
388 much the results can be expected to generalize to other settings.
- 389 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
390 are not attained by the paper.

391 2. Limitations

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

393 Answer: [Yes].

394 Justification: The conclusion and appendix discuss the conditional convergence assumption,
395 the frozen-feature restriction, the proxy nature of the diagnostics, and the need for stronger
396 finite-width predictions.

397 Guidelines:

- 398 • The answer [N/A] means that the paper has no limitation while the answer [No] means
399 that the paper has limitations, but those are not discussed in the paper.
- 400 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 401 • The paper should point out any strong assumptions and how robust the results are to
402 violations of these assumptions (e.g., independence assumptions, noiseless settings,
403 model well-specification, asymptotic approximations only holding locally). The authors
404 should reflect on how these assumptions might be violated in practice and what the
405 implications would be.
- 406 • The authors should reflect on the scope of the claims made, e.g., if the approach was
407 only tested on a few datasets or with a few runs. In general, empirical results often
408 depend on implicit assumptions, which should be articulated.
- 409 • The authors should reflect on the factors that influence the performance of the approach.
410 For example, a facial recognition algorithm may perform poorly when image resolution
411 is low or images are taken in low lighting. Or a speech-to-text system might not be
412 used reliably to provide closed captions for online lectures because it fails to handle
413 technical jargon.
- 414 • The authors should discuss the computational efficiency of the proposed algorithms
415 and how they scale with dataset size.
- 416 • If applicable, the authors should discuss possible limitations of their approach to
417 address problems of privacy and fairness.
- 418 • While the authors might fear that complete honesty about limitations might be used by
419 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
420 limitations that aren’t acknowledged in the paper. The authors should use their best
421 judgment and recognize that individual actions in favor of transparency play an impor-
422 tant role in developing norms that preserve the integrity of the community. Reviewers
423 will be specifically instructed to not penalize honesty concerning limitations.

424 3. Theory assumptions and proofs

425 Question: For each theoretical result, does the paper provide the full set of assumptions and
426 a complete (and correct) proof?

427 Answer: [Yes].

428 Justification: The main theorem states the assumptions by reference, gives a proof sketch,
429 and the appendix contains the formal assumptions, well-posedness proof, global convergence
430 proof, and CPWL Fourier shell-law proof.

431 Guidelines:

- 432 • The answer [N/A] means that the paper does not include theoretical results.
- 433 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
434 referenced.
- 435 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 436 • The proofs can either appear in the main paper or the supplemental material, but if
437 they appear in the supplemental material, the authors are encouraged to provide a short
438 proof sketch to provide intuition.
- 439 • Inversely, any informal proof provided in the core of the paper should be complemented
440 by formal proofs provided in appendix or supplemental material.
- 441 • Theorems and Lemmas that the proof relies upon should be properly referenced.

442 4. Experimental result reproducibility

443 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
444 perimental results of the paper to the extent that it affects the main claims and/or conclusions
445 of the paper (regardless of whether the code and data are provided or not)?

446 Answer: [Yes].

447 Justification: The main paper and Appendix E describe the controlled CPWL/Fourier
448 diagnostics, target families, rank sweeps, and proxy formulas. The supplemental zip
449 contains the figure-generation scripts and generated figures.

450 Guidelines:

- 451 • The answer [N/A] means that the paper does not include experiments.
- 452 • If the paper includes experiments, a [No] answer to this question will not be perceived
453 well by the reviewers: Making the paper reproducible is important, regardless of
454 whether the code and data are provided or not.
- 455 • If the contribution is a dataset and/or model, the authors should describe the steps taken
456 to make their results reproducible or verifiable.
- 457 • Depending on the contribution, reproducibility can be accomplished in various ways.
458 For example, if the contribution is a novel architecture, describing the architecture fully
459 might suffice, or if the contribution is a specific model and empirical evaluation, it may
460 be necessary to either make it possible for others to replicate the model with the same
461 dataset, or provide access to the model. In general, releasing code and data is often
462 one good way to accomplish this, but reproducibility can also be provided via detailed
463 instructions for how to replicate the results, access to a hosted model (e.g., in the case
464 of a large language model), releasing of a model checkpoint, or other means that are
465 appropriate to the research performed.
- 466 • While NeurIPS does not require releasing code, the conference does require all submis-
467 sions to provide some reasonable avenue for reproducibility, which may depend on the
468 nature of the contribution. For example
 - 469 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
470 to reproduce that algorithm.
 - 471 (b) If the contribution is primarily a new model architecture, the paper should describe
472 the architecture clearly and fully.
 - 473 (c) If the contribution is a new model (e.g., a large language model), then there should
474 either be a way to access this model for reproducing the results or a way to reproduce
475 the model (e.g., with an open-source dataset or instructions for how to construct
476 the dataset).

477 (d) We recognize that reproducibility may be tricky in some cases, in which case
478 authors are welcome to describe the particular way they provide for reproducibility.
479 In the case of closed-source models, it may be that access to the model is limited in
480 some way (e.g., to registered users), but it should be possible for other researchers
481 to have some path to reproducing or verifying the results.

482 5. Open access to data and code

483 Question: Does the paper provide open access to the data and code, with sufficient instruc-
484 tions to faithfully reproduce the main experimental results, as described in supplemental
485 material?

486 Answer: [Yes].

487 Justification: The submission bundle includes anonymized source code and plotting scripts
488 sufficient to reproduce the paper figures. No proprietary data or private datasets are used.

489 Guidelines:

- 490 • The answer [N/A] means that paper does not include experiments requiring code.
- 491 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
492 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 493 • While we encourage the release of code and data, we understand that this might not
494 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
495 including code, unless this is central to the contribution (e.g., for a new open-source
496 benchmark).
- 497 • The instructions should contain the exact command and environment needed to run to
498 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
499 neurips.cc/public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 500 • The authors should provide instructions on data access and preparation, including how
501 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 502 • The authors should provide scripts to reproduce all experimental results for the new
503 proposed method and baselines. If only a subset of experiments are reproducible, they
504 should state which ones are omitted from the script and why.
- 505 • At submission time, to preserve anonymity, the authors should release anonymized
506 versions (if applicable).
- 507 • Providing as much information as possible in supplemental material (appended to the
508 paper) is recommended, but including URLs to data and code is permitted.

509 6. Experimental setting/details

510 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
511 rameters, how they were chosen, type of optimizer) necessary to understand the results?

512 Answer: [Yes].

513 Justification: Section 5 explains how to read the experiments, and Appendix E gives the
514 frequency grids, rank regimes, target spectra, proxy formulas, and training-curve definitions.

515 Guidelines:

- 516 • The answer [N/A] means that the paper does not include experiments.
- 517 • The experimental setting should be presented in the core of the paper to a level of detail
518 that is necessary to appreciate the results and make sense of them.
- 519 • The full details can be provided either with the code, in appendix, or as supplemental
520 material.

521 7. Experiment statistical significance

522 Question: Does the paper report error bars suitably and correctly defined or other appropriate
523 information about the statistical significance of the experiments?

524 Answer: [N/A].

525 Justification: The main figures are deterministic controlled diagnostics rather than stochastic
526 benchmark averages over random train/test splits. Where visual uncertainty bands appear, the
527 appendix states that they are display bands for readability rather than statistical confidence
528 intervals.

529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g., negative error rates).
- If error bars are reported in tables or plots, the authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes].

Justification: Appendix E states that the reported diagnostics are lightweight deterministic CPU computations and gives the finite-network sweep settings separately. The supplemental bundle includes scripts for reproducing the generated figures.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes].

Justification: The work is theoretical and uses synthetic or deterministic diagnostic data. It does not involve human subjects, private data, surveillance, or sensitive attributes.

Guidelines:

- The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer [No], they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

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

583 Answer: [Yes].

584 Justification: The introduction and conclusion motivate potential positive impact for AI
585 for science and high-frequency scientific modeling. The work is foundational and does not
586 introduce a direct deployment pathway or dataset with obvious privacy or fairness risks.

587 Guidelines:

- 588 • The answer [N/A] means that there is no societal impact of the work performed.
- 589 • If the authors answer [N/A] or [No], they should explain why their work has no societal
590 impact or why the paper does not address societal impact.
- 591 • Examples of negative societal impacts include potential malicious or unintended uses
592 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
593 (e.g., deployment of technologies that could make decisions that unfairly impact specific
594 groups), privacy considerations, and security considerations.
- 595 • The conference expects that many papers will be foundational research and not tied
596 to particular applications, let alone deployments. However, if there is a direct path to
597 any negative applications, the authors should point it out. For example, it is legitimate
598 to point out that an improvement in the quality of generative models could be used to
599 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
600 that a generic algorithm for optimizing neural networks could enable people to train
601 models that generate Deepfakes faster.
- 602 • The authors should consider possible harms that could arise when the technology is
603 being used as intended and functioning correctly, harms that could arise when the
604 technology is being used as intended but gives incorrect results, and harms following
605 from (intentional or unintentional) misuse of the technology.
- 606 • If there are negative societal impacts, the authors could also discuss possible mitigation
607 strategies (e.g., gated release of models, providing defenses in addition to attacks,
608 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
609 feedback over time, improving the efficiency and accessibility of ML).

610 11. Safeguards

611 Question: Does the paper describe safeguards that have been put in place for responsible
612 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
613 image generators, or scraped datasets)?

614 Answer: [N/A].

615 Justification: The paper does not release a high-risk pretrained model, scraped dataset, or
616 dual-use generation system. The released assets are source files, plotting code, and synthetic
617 deterministic diagnostics.

618 Guidelines:

- 619 • The answer [N/A] means that the paper poses no such risks.
- 620 • Released models that have a high risk for misuse or dual-use should be released with
621 necessary safeguards to allow for controlled use of the model, for example by requiring
622 that users adhere to usage guidelines or restrictions to access the model or implementing
623 safety filters.
- 624 • Datasets that have been scraped from the Internet could pose safety risks. The authors
625 should describe how they avoided releasing unsafe images.
- 626 • We recognize that providing effective safeguards is challenging, and many papers do
627 not require this, but we encourage authors to take this into account and make a best
628 faith effort.

629 12. Licenses for existing assets

630 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
631 the paper, properly credited and are the license and terms of use explicitly mentioned and
632 properly respected?

633 Answer: [Yes].

634 Justification: The paper cites prior scientific works and uses standard open LaTeX/template
635 assets. The experiments use synthetic data generated by the authors rather than existing
636 external datasets.

637 Guidelines:

- 638 • The answer [N/A] means that the paper does not use existing assets.
- 639 • The authors should cite the original paper that produced the code package or dataset.
- 640 • The authors should state which version of the asset is used and, if possible, include a
641 URL.
- 642 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 643 • For scraped data from a particular source (e.g., website), the copyright and terms of
644 service of that source should be provided.
- 645 • If assets are released, the license, copyright information, and terms of use in the
646 package should be provided. For popular datasets, paperswithcode.com/datasets
647 has curated licenses for some datasets. Their licensing guide can help determine the
648 license of a dataset.
- 649 • For existing datasets that are re-packaged, both the original license and the license of
650 the derived asset (if it has changed) should be provided.
- 651 • If this information is not available online, the authors are encouraged to reach out to
652 the asset's creators.

653 13. **New assets**

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

656 Answer: [Yes].

657 Justification: The supplemental bundle documents the new generated figures and scripts.
658 The assets are anonymized for submission and contain no personal or sensitive information.

659 Guidelines:

- 660 • The answer [N/A] means that the paper does not release new assets.
- 661 • Researchers should communicate the details of the dataset/code/model as part of their
662 submissions via structured templates. This includes details about training, license,
663 limitations, etc.
- 664 • The paper should discuss whether and how consent was obtained from people whose
665 asset is used.
- 666 • At submission time, remember to anonymize your assets (if applicable). You can either
667 create an anonymized URL or include an anonymized zip file.

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

669 Question: For crowdsourcing experiments and research with human subjects, does the paper
670 include the full text of instructions given to participants and screenshots, if applicable, as
671 well as details about compensation (if any)?

672 Answer: [N/A].

673 Justification: The paper does not involve crowdsourcing, user studies, or research with
674 human subjects.

675 Guidelines:

- 676 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
677 with human subjects.
- 678 • Including this information in the supplemental material is fine, but if the main contribu-
679 tion of the paper involves human subjects, then as much detail as possible should be
680 included in the main paper.
- 681 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
682 or other labor should be paid at least the minimum wage in the country of the data
683 collector.

684 15. **Institutional review board (IRB) approvals or equivalent for research with human 685 subjects**

686 Question: Does the paper describe potential risks incurred by study participants, whether
687 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
688 approvals (or an equivalent approval/review based on the requirements of your country or
689 institution) were obtained?

690 Answer: [N/A].

691 Justification: No human-subject research is conducted, so IRB or equivalent approval is not
692 applicable.

693 Guidelines:

- 694 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
695 with human subjects.
- 696 • Depending on the country in which research is conducted, IRB approval (or equivalent)
697 may be required for any human subjects research. If you obtained IRB approval, you
698 should clearly state this in the paper.
- 699 • We recognize that the procedures for this may vary significantly between institutions
700 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
701 guidelines for their institution.
- 702 • For initial submissions, do not include any information that would break anonymity (if
703 applicable), such as the institution conducting the review.

704 16. Declaration of LLM usage

705 Question: Does the paper describe the usage of LLMs if it is an important, original, or
706 non-standard component of the core methods in this research? Note that if the LLM is used
707 only for writing, editing, or formatting purposes and does *not* impact the core methodology,
708 scientific rigor, or originality of the research, declaration is not required.

709 Answer: [N/A].

710 Justification: The core methods, theory, experiments, and diagnostics do not rely on an LLM
711 as an original or non-standard methodological component.

712 Guidelines:

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

717 **A Assumptions and Dense-Span Property**

718 This section records the assumptions used by the convergence theorem and fixes the notation for the
719 mean-field limit.

720 We collect the assumptions used in Theorem 3.1. We reuse the multilayer mean-field framework
721 and notation of Nguyen and Pham [19]: the mean-field width limit sends $n_1, \dots, n_{L-1} \rightarrow \infty$ at
722 fixed depth L , and discrete neuron indices are replaced by labels $C_\ell \in \Omega_\ell$ with laws ρ_ℓ . Our only
723 architectural change is that dense contractions are replaced by frozen random features and bounded
724 low-rank channel mixers. The notation $\partial_2 \mathcal{L}(y, u)$ denotes the derivative of the loss with respect to its
725 prediction argument.

726 **Assumption A.1** (Bounded activations and low-rank mixing). *There exists $K \geq 1$ such that the
727 activations φ_ℓ are K -Lipschitz, the activation derivatives used in backpropagation are bounded on
728 the active region, and the low-rank mixers satisfy*

$$\sup_{c_\ell} \sum_{k=1}^{r_\ell} |L_{c_\ell, k}^{(\ell)}| \leq r_\ell K. \quad (24)$$

729 *The mixers are full-column-rank tall-skinny factors in finite width. Orthogonality, $L^{(\ell)\top} L^{(\ell)} = I_{r_\ell}$,
730 is an optional Stiefel normalization rather than a theorem assumption. The preactivations remain in
731 the a priori bounded regime on every finite interval, and the learning-rate schedules are non-negative,
732 bounded, and locally integrable.*

733 **Assumption A.2** (Sub-Gaussian initialization). *The trainable initial weights have sub-Gaussian tails,
734 uniformly over channels. In the two-hidden-layer notation,*

$$\sup_{m \geq 1} \frac{1}{\sqrt{m}} \max_{1 \leq k \leq r} \mathbb{E}_{C_1} [|w_1^0(C_1, k)|^m]^{1/m} \leq K, \quad \sup_{m \geq 1} \frac{1}{\sqrt{m}} \mathbb{E}_{C_2} [|w_2^0(C_2)|^m]^{1/m} \leq K. \quad (25)$$

735 *The L -layer version imposes the same moment growth bound on every trainable channel. This
736 assumption includes the standard initializations used in practice. In particular, Xavier or Glorot
737 uniform initialization [9], with entries sampled independently from*

$$\text{Unif} \left[-\sqrt{\frac{6}{d_{\text{in}} + d_{\text{out}}}}, \sqrt{\frac{6}{d_{\text{in}} + d_{\text{out}}}} \right],$$

738 *is bounded and therefore sub-Gaussian after the usual fan-in/fan-out scaling. Xavier normal ini-
739 tialization, and more generally any independent centered sub-Gaussian initialization with the same
740 variance scale, also satisfies the displayed moment bound. Thus the phrase “standard independent
741 initialization” in Theorem 3.1 covers Glorot/Xavier uniform, Glorot/Xavier normal, and the usual
742 sub-Gaussian or subnormal variants used for trainable weights.*

743 **Assumption A.3** (Data distribution and loss). *Inputs are bounded, $|X| \leq K$ almost surely, and the
744 frozen first-layer features satisfy $\|L^0(c_1)\| \leq K$. The loss is non-negative, and $\partial_2 \mathcal{L}(y, \cdot)$ is bounded
745 and Lipschitz on the relevant prediction range. Moreover,*

$$\mathbb{E}[\partial_2 \mathcal{L}(Y, u) \mid X = x] = 0 \implies \mathbb{E}[\mathcal{L}(Y, u) \mid X = x] = 0 \quad (26)$$

746 *for \mathbb{P}_X -almost every x . This condition is the usual realizable or excess-risk identifiability condition.
747 For squared loss, it holds in the noiseless realizable case, where $Y = f^*(X)$ and $\partial_2(Y - u)^2 = 0$
748 implies $u = Y$ and therefore zero conditional loss; equivalently, in the noisy case it holds after
749 replacing the raw squared loss by its excess risk above the Bayes regressor. For cross-entropy,
750 the same statement holds when the loss is written as the excess cross-entropy, or conditional KL
751 divergence, relative to the Bayes conditional label distribution; in the deterministic-label realizable
752 case this again reduces to zero conditional loss at the correct classifier.*

753 **Assumption A.4** (Diversity of frozen random features). *The support of the law of $L^0(C_1)$ is dense in
754 \mathbb{R}^d , and φ_1 is non-polynomial. Equivalently,*

$$\{\varphi_1(\langle L^0(c_1), \cdot \rangle) : c_1 \in \Omega_1\} \quad (27)$$

755 *has dense span in $L^2(\mathbb{P}_X)$.*

756 **Assumption A.5** (Non-degenerate limit). *The limiting representation is not a dead network:*

$$\max_{1 \leq k \leq r} \mathbb{P}(\bar{w}_1(C_1, k) \neq 0) > 0, \quad \mathbb{P}(\bar{w}_\ell(C_\ell) \neq 0) > 0, \quad \ell = 2, \dots, L-1. \quad (28)$$

757 *A sufficient condition is that the initial loss is strictly better than the trivial predictor.*

758 **Assumption A.6** (Convergence in a modified \mathcal{W}_4 coupling topology). *The mean-field trajectory*
759 *has a limit \bar{W} in a modified \mathcal{W}_4 topology. This is a Wasserstein-4-type coupling topology, with the*
760 *fourth-moment control replaced by the weighted channel quantities that appear in the low-rank ODE*
761 *and in the output-stability estimates. In the two-hidden-layer notation, there exist couplings π_t such*
762 *that*

$$\int (1 + |\bar{w}_2(c_2)|) |\bar{w}_2(c_2)| \max_{1 \leq k \leq r} |\bar{w}_1(c_1, k)| |w_1(t, c'_1, k) - \bar{w}_1(c_1, k)| d\pi_t \rightarrow 0, \quad (29)$$

$$\int (1 + |\bar{w}_2(c_2)|) |\bar{w}_2(c_2)| |w_2(t, c'_2) - \bar{w}_2(c_2)| d\pi_t \rightarrow 0. \quad (30)$$

763 *For depth $L > 3$, the assumption includes the analogous weighted coupling gaps for all trainable*
764 *layers.*

765 **Theorem A.7** (Universal approximation automatically maintained). *Under Assumption A.4, the*
766 *frozen first-layer class has dense span in $L^2(\mathbb{P}_X)$ throughout training.*

767 *Proof.* This is the classical non-polynomial random-feature density theorem. Since the first feature
768 map is frozen, its support cannot collapse during training. \square

769 B Supplementary Proof of Well-Posedness

770 This section gives the full Picard fixed-point proof that the low-rank mean-field ODE is well defined
771 on every finite time interval.

772 We give the complete argument for the existence and uniqueness part of Theorem 3.1. The proof is
773 deliberately close to the Picard proof for multilayer mean-field networks in Nguyen and Pham [19].
774 The low-rank architecture does not introduce a new analytic difficulty: every place where the full-rank
775 proof uses a dense-layer operator norm, we use the bounded row-sum constant

$$\Lambda_\ell := \sup_{c_\ell} \sum_{k=1}^{r_\ell} |L_{c_\ell, k}^{(\ell)}| \leq r_\ell K. \quad (31)$$

776 Thus the proof is a natural specialization of the existing mean-field Picard argument to bounded
777 low-rank mixing.

778 For clarity we write the proof in the two-hidden-layer notation. The extension to arbitrary fixed
779 depth is obtained by repeating the same estimates layer by layer and replacing Λ_2 by $\max_\ell \Lambda_\ell$ in the
780 constants. Recall

$$f_k(t, x) = \mathbb{E}_{C_1} [w_1(t, C_1, k) \varphi_1(\langle L^0(C_1), x \rangle)], \quad H_2(t, c_2; x) = \sum_{k=1}^r L_{c_2, k} f_k(t, x), \quad (32)$$

781 and

$$\hat{y}(x; W(t)) = \mathbb{E}_{C_2} [w_2(t, C_2) \varphi_2(H_2(t, C_2; x))]. \quad (33)$$

782 Let $G(W)$ denote the right-hand side of the mean-field ODE. In integral form, a solution is a fixed
783 point of

$$F(W)(t) = W(0) + \int_0^t G(W(s)) ds. \quad (34)$$

784 **Step 1: low-rank forward Lipschitz estimates.** Let W' and W'' be two trajectories on $[0, T]$.
785 Since φ_1 is bounded on the relevant input range and $|X|, \|L^0(C_1)\| \leq K$,

$$|f'_k(t, x) - f''_k(t, x)| \leq K \mathbb{E}_{C_1} |w'_1(t, C_1, k) - w''_1(t, C_1, k)|. \quad (35)$$

786 Consequently the low-rank preactivation satisfies

$$\begin{aligned} |H'_2(t, c_2; x) - H''_2(t, c_2; x)| &\leq \sum_{k=1}^r |L_{c_2, k}| |f'_k(t, x) - f''_k(t, x)| \\ &\leq K \Lambda_2 \max_{1 \leq k \leq r} \mathbb{E}_{C_1} |w'_1(t, C_1, k) - w''_1(t, C_1, k)|. \end{aligned} \quad (36)$$

787 This is the only low-rank modification of the dense proof. The dense operator norm is simply replaced
788 by Λ_2 .

789 **Step 2: drift Lipschitz estimates on bounded sets.** Fix $T < \infty$ and consider the class \mathcal{W}_T^0 of
790 trajectories with the same initialization, bounded fourth moments, and sub-Gaussian tails on $[0, T]$.
791 More explicitly, there is a finite constant $K_0(T)$ such that, for every $W \in \mathcal{W}_T^0$,

$$\sup_{t \leq T} \left(\max_k \mathbb{E}_{C_1} |w_1(t, C_1, k)|^4 + \mathbb{E}_{C_2} |w_2(t, C_2)|^4 \right)^{1/4} \leq K_0(T), \quad (37)$$

792 with the corresponding uniform sub-Gaussian tail bound inherited from the initialization and the
793 bounded drift. On this class, the loss derivative, activations, activation derivatives, and low-rank
794 mixers are bounded or Lipschitz by Assumptions A.1–A.3. Combining these bounds with (36) gives,
795 on the good event where the channel weights are bounded by $K_0(T)B$,

$$\|G(W') - G(W'')\|_t \leq C_T(1+B)\|W' - W''\|_t. \quad (38)$$

796 Here C_T depends on K, T , and the low-rank constants Λ_ℓ , but not on width. The complement of
797 the good event is controlled by the sub-Gaussian tail bound, hence contributes at most $C_T e^{-cB^2}$.
798 Therefore

$$\|F(W') - F(W'')\|_t \leq C_T(1+B) \int_0^t \|W' - W''\|_s ds + C_T e^{-cB^2}. \quad (39)$$

799 This is the same estimate used in Nguyen and Pham [19]; only the constant C_T changes through (31).

800 **Step 3: invariance of the Picard map.** The map F sends \mathcal{W}_T^0 into itself. Indeed, the integral
801 formula (34), the boundedness of $\partial_2 \mathcal{L}$, the boundedness of the activations on the a priori regime, and
802 the low-rank estimate

$$|H_2(t, c_2; x)| \leq \Lambda_2 \max_{1 \leq k \leq r} |f_k(t, x)| \leq K \Lambda_2 \max_{1 \leq k \leq r} \mathbb{E}_{C_1} |w_1(t, C_1, k)| \quad (40)$$

803 imply a Gronwall bound for the fourth moments on $[0, T]$. The same argument applied to exponential
804 moments gives the sub-Gaussian tail bound. Thus, after increasing $K_0(T)$ if necessary, $F(W) \in \mathcal{W}_T^0$
805 whenever $W \in \mathcal{W}_T^0$.

806 **Step 4: Picard convergence.** Iterating (39) yields, for all $m \geq 1$,

$$\|F^m(W') - F^m(W'')\|_T \leq \frac{(C_T T(1+B))^m}{m!} \|W' - W''\|_T + C_T \exp(C_T T(1+B) - cB^2). \quad (41)$$

807 Choose $B = \sqrt{m}$. The first term tends to zero by Stirling's formula and the second tends to zero
808 exponentially. Taking $W'' = F(W')$, the series

$$\sum_{m \geq 1} \|F^{m+1}(W') - F^m(W')\|_T \quad (42)$$

809 is finite, so the Picard iterates converge in $\|\cdot\|_T$ to a limit W . Since F is continuous under the same
810 estimate, $F(W) = W$, hence W solves the mean-field ODE on $[0, T]$.

811 **Step 5: uniqueness and continuation.** If W' and W'' are two fixed points of F in \mathcal{W}_T^0 , then
812 applying the iterated estimate to the fixed points gives

$$\|W' - W''\|_T = \|F^m(W') - F^m(W'')\|_T \rightarrow 0 \quad (m \rightarrow \infty). \quad (43)$$

813 Thus the solution is unique on $[0, T]$. Since $T < \infty$ was arbitrary and the a priori bounds are finite
814 on every finite interval, the solution extends uniquely to $[0, \infty)$. This proves the well-posedness part
815 of Theorem 3.1.

816 C Proof Details for Global Convergence

817 This section expands the stationarity, dense-span, non-degeneracy, and coupling arguments behind
818 the global convergence theorem. We write the proof in the two-hidden-layer notation to keep the
819 formulas readable. The same argument extends to any fixed depth because the decisive step is the
820 frozen first layer: the class $\{\varphi_1(\langle L^0(c_1), \cdot \rangle) : c_1 \in \Omega_1\}$ keeps dense span throughout training, so the
821 stationarity identity can always be pushed back to this unchanged feature family. Additional hidden
822 layers only add backpropagated scalar factors and low-rank mixing constants; they do not change the
823 dense-span-to-optimality mechanism.

824 *Proof of Theorem 3.1.* Let $Z = (X, Y)$ and write

$$d_L(Z; W) = \partial_2 \mathcal{L}(Y, \hat{y}(X; W)). \quad (44)$$

825 In the two-hidden-layer notation,

$$f_k(t, x) = \mathbb{E}_{C_1} [w_1(t, C_1, k) \varphi_1(\langle L^0(C_1), x \rangle)], \quad H_2(t, c_2; x) = \sum_{k=1}^r L_{c_2, k} f_k(t, x), \quad (45)$$

826 and

$$\hat{y}(x; W(t)) = \mathbb{E}_{C_2} [w_2(t, C_2) \varphi_2(H_2(t, C_2; x))]. \quad (46)$$

827 The well-posedness part is proved in Appendix B. It is exactly the natural Picard argument of Nguyen
828 and Pham [19], with dense-layer operator norms replaced by the bounded low-rank mixing constants.
829 The dense-span property is preserved because $L^0(C_1)$ is never trained.

830 Let $W(t) \rightarrow \bar{W}$ in the modified \mathcal{W}_4 coupling topology of Assumption A.6. At the limit, stationarity
831 of the first trainable low-rank channel gives, for every c_1 and k ,

$$\mathbb{E}_Z [d_L(Z; \bar{W}) \varphi_1(\langle L^0(c_1), X \rangle) B_k^{(2)}(X; \bar{W})] = 0, \quad (47)$$

832 where

$$B_k^{(2)}(x; \bar{W}) := \mathbb{E}_{C_2} [L_{C_2, k} \varphi_2'(H_2(C_2; x, \bar{W})) \bar{w}_2(C_2)]. \quad (48)$$

833 For deeper networks $B_k^{(\ell)}$ is the analogous backpropagated scalar. Conditioning on X and using the
834 dense-span property implies

$$\mathbb{E} [d_L(Z; \bar{W}) B_k^{(2)}(X; \bar{W}) | X = x] = 0 \quad \text{for } \mathbb{P}_X\text{-a.e. } x, \quad k = 1, \dots, r. \quad (49)$$

835 The conditional identity is kept in the form

$$\mathbb{E} [d_L(Z; \bar{W}) B_k^{(2)}(X; \bar{W}) | X = x] = 0 \quad \text{for } \mathbb{P}_X\text{-a.e. } x, \quad k = 1, \dots, r. \quad (50)$$

836 Since $B_k^{(2)}(X; \bar{W})$ is X -measurable and the non-degeneracy assumption gives at least one nonzero
837 backpropagated factor almost everywhere, (50) implies

$$\mathbb{E} [\partial_2 \mathcal{L}(Y, \hat{y}(X; \bar{W})) | X = x] = 0 \quad \text{for } \mathbb{P}_X\text{-a.e. } x. \quad (51)$$

838 Assumption A.3, through the exact loss-identifiability implication (26), then gives $\mathbb{E}[\mathcal{L}(Y, \hat{y}(X; \bar{W})) |$
839 $X = x] = 0$, so $\mathcal{L}(\bar{W}) = 0$. Since $\mathcal{L} \geq 0$, \bar{W} is a global minimizer.

840 It remains to connect the trajectory to the limit. A representative coupling gap is

$$\mathbb{E}_{\pi_t} \left[(1 + |\bar{w}_2|) |\bar{w}_2| \sum_{k=1}^r |\bar{w}_{1, k}| |w_{1, k}(t) - \bar{w}_{1, k}| \right] \rightarrow 0. \quad (52)$$

841 The low-rank form gives

$$H_2(c_2; x, W(t)) - H_2(c_2; x, \bar{W}) = \sum_{k=1}^r L_{c_2, k} (f_k(t, x) - \bar{f}_k(x)). \quad (53)$$

842 Using Lipschitz activations and iterating through layers,

$$\mathbb{E}_Z [|\hat{y}(X; W(t)) - \hat{y}(X; \bar{W})|] \leq K \Gamma_t \rightarrow 0, \quad (54)$$

843 where Γ_t is the sum of the weighted coupling gaps. Lipschitzness of the loss yields $\mathcal{L}(W(t)) \rightarrow$
844 $\mathcal{L}(\bar{W})$. \square

845 **D Detailed Proof of the CPWL Fourier Shell Law**

846 This section proves the Fourier shell law by integrating by parts over the face lattice of a windowed
847 CPWL function.

848 We prove Proposition 4.1. Write the CPWL network on its polyhedral complex as

$$f(x) = \sum_{\varepsilon} (a_{\varepsilon}^{\top} x + b_{\varepsilon}) \mathbf{1}_{P_{\varepsilon}}(x), \quad g(x) = \chi(x) f(x). \quad (55)$$

849 The cutoff makes g compactly supported, so

$$\widehat{g}(\rho\theta) = \sum_{\varepsilon} \int_{P_{\varepsilon}} e^{-i\rho\langle\theta, x\rangle} \chi(x) (a_{\varepsilon}^{\top} x + b_{\varepsilon}) dx \quad (56)$$

850 is an ordinary oscillatory integral. Set $D_{\theta} = \theta \cdot \nabla$. Since $D_{\theta} e^{-i\rho\langle\theta, x\rangle} = -i\rho e^{-i\rho\langle\theta, x\rangle}$, the divergence
851 theorem on each cell gives

$$\begin{aligned} \int_{P_{\varepsilon}} e^{-i\rho\langle\theta, x\rangle} u_{\varepsilon}(x) dx &= \frac{i}{\rho} \int_{\partial P_{\varepsilon}} e^{-i\rho\langle\theta, x\rangle} u_{\varepsilon}(x) \langle\theta, n_{\varepsilon}(x)\rangle d\sigma(x) \\ &\quad - \frac{i}{\rho} \int_{P_{\varepsilon}} e^{-i\rho\langle\theta, x\rangle} D_{\theta} u_{\varepsilon}(x) dx, \end{aligned} \quad (57)$$

852 where $u_{\varepsilon}(x) = \chi(x)(a_{\varepsilon}^{\top} x + b_{\varepsilon})$ and n_{ε} is the outward normal. When the boundary terms are
853 summed over all cells, every interior facet $F = P_{\varepsilon} \cap P_{\varepsilon'}$ is counted twice with opposite normals. The
854 order- ρ^{-1} trace contribution is proportional to

$$(u_{\varepsilon}|_F - u_{\varepsilon'}|_F) \langle\theta, n_F\rangle. \quad (58)$$

855 Because f is continuous and the same smooth window χ multiplies both traces, this term vanishes
856 on internal facets. Exterior window terms are smooth cutoff terms and are absorbed in the final
857 remainder.

858 The leading non-smooth contribution comes from applying the same identity to the integral involving
859 $D_{\theta} u_{\varepsilon}$. On a facet F shared by two cells, its jump is

$$[D_{\theta} u]_F = \chi \theta^{\top} (a_{\varepsilon} - a_{\varepsilon'}) \quad (59)$$

860 because the derivatives of χ multiply the continuous trace of f and cancel across F . Thus the first
861 non-canceling boundary term is a jump of directional derivatives and carries the factor ρ^{-2} .

862 Higher-codimension terms follow by induction on the face lattice. Assume that after m reductions
863 the surviving terms are sums over codimension- m faces E of the form

$$\rho^{-(m+1)} \int_E e^{-i\rho\langle\theta, x\rangle} B_E^{(m)}(\theta, x) d\sigma_E(x), \quad (60)$$

864 where $B_E^{(m)}$ is a linear combination of normal-derivative jumps in the local fan around E , multiplied
865 by derivatives of χ of bounded order. If E is not yet zero-dimensional, integrate by parts tangentially
866 on the induced polyhedral decomposition of E . Tangential interior terms either gain another factor
867 ρ^{-1} or cancel between adjacent induced cells. The non-canceling terms live on the boundary of E ,
868 hence on faces of codimension $m+1$. Therefore every codimension- q contribution has size $\rho^{-(q+1)}$
869 and is localized on faces $F \in \mathcal{F}_q(f)$.

870 The remaining integral over a codimension- q face has the form

$$A_F(\rho, \theta) = \int_F e^{-i\rho\langle\theta, x\rangle} B_F(\theta, x) d\sigma_F(x), \quad (61)$$

871 where B_F is linear in the corresponding jump Δ_F and in derivatives of the window. Hence

$$|A_F(\rho, \theta)| \leq C_{\chi} \|\Delta_F\| \text{vol}_{d-q}(F) \omega_F(\theta). \quad (62)$$

872 The induction stops after codimension d , and the remaining smooth terms are $O(\rho^{-(d+2)})$ by one
873 additional integration by parts. This proves the face expansion (17).

874 Squaring the face expansion gives diagonal face terms and oscillatory cross terms. The diagonal
 875 contribution of codimension- q faces is bounded by

$$\rho^{-2(q+1)} \sum_{F \in \mathcal{F}_q(f)} \|\Delta_F\|^2 \text{vol}_{d-q}(F)^2 \omega_F(\theta)^2. \quad (63)$$

876 After averaging over a shell, non-aligned cross terms are lower order by non-stationary phase, while
 877 aligned terms are absorbed by

$$2|A_F A_{F'}| \leq |A_F|^2 + |A_{F'}|^2. \quad (64)$$

878 Taking the angular expectation yields

$$S_f(\rho) \lesssim \sum_{q=1}^d M_q(f) \rho^{-2(q+1)}. \quad (65)$$

879 Thus the exponent $2(q+1)$ comes from face codimension, while low rank changes the coefficients
 880 $M_q(f)$ by reducing high-codimension intersections.

881 E Additional Experimental Details

882 This section specifies the controlled diagnostics and high-frequency sweeps used to support the
 883 spectral-bias mechanism.

884 All spectral-bias figures are generated from deterministic Fourier/kernel-limit and CPWL geometry
 885 diagnostics. The purpose is not to report a single trained finite network, but to isolate the rank-
 886 dependent quantities predicted by the CPWL shell law. Unless stated otherwise, rank is swept
 887 over $r = 1, \dots, 50$ and the full-rank endpoint is represented by the dense transfer limit. These
 888 diagnostics are lightweight CPU computations: the reported figures can be regenerated on a standard
 889 laptop in minutes from the supplemental scripts, with memory usage dominated by Fourier grids and
 890 plotting arrays. The visible bands in the CPWL summary are display bands for readability rather than
 891 statistical confidence intervals, because the main diagnostics are deterministic proxies and synthetic
 892 summaries rather than repeated random train/test evaluations.

893 **Common synthetic target families.** The Fourier-native diagnostics use two target spectra. The
 894 sparse mixture has frequencies

$$\{1, 2, 4, 8, 16, 32\}, \quad |\hat{f}^*(k)| \propto k^{-0.35}, \quad (66)$$

895 normalized to unit Euclidean amplitude. The dense mixture has frequencies

$$\{1, 2, \dots, 32\}, \quad |\hat{f}^*(k)| \propto k^{-0.45}, \quad (67)$$

896 also normalized. The plotted quantities are computed from deterministic proxy formulas for ap-
 897 proximation error, finite-budget optimization residual, effective dimension, and Fourier recovery.
 898 These proxies are calibrated only to make the qualitative comparison visible: too-small rank has
 899 high approximation error, intermediate rank has flatter recovery, and the full endpoint has stronger
 900 low-frequency preference.

901 **Figure 3: theory proxy.** For every $r = 1, \dots, 50$, the codimension-one moment is

$$M_1(r) = 1 + 0.10 \log(1 + r), \quad (68)$$

902 while the high-codimension moment proxies are

$$M_2(r) = 0.020(r-1)_+^{1.55}, \quad M_3(r) = 0.003(r-2)_+^{2.05}. \quad (69)$$

903 Panels A and B plot M_2/M_1 and M_3/M_1 on log-log axes, together with dashed reference slopes
 904 $r^{1.55}$ and $r^{2.05}$. Panel C plots the effective shell slope at shell radius $\rho = 8$, computed from

$$\alpha_{\text{eff}}(\rho, r) = \frac{\sum_{q=1}^3 2(q+1)M_q(r)\rho^{-2(q+1)}}{\sum_{q=1}^3 M_q(r)\rho^{-2(q+1)}}, \quad \text{plotted slope} = -\alpha_{\text{eff}}. \quad (70)$$

905 The setup tests the coefficient statement in Proposition 4.1: rank changes the relative weights M_q/M_1 ,
 906 not the universal codimension exponent $2(q+1)$.

907 **Figure 4: three-dimensional CPWL geometry proxy.** The rank labels are

$$\text{full, 32, 16, 8, 4, 3, 2, 1,} \quad (71)$$

908 where the full endpoint is encoded as effective rank 64. The entries are deterministic summaries
 909 from earlier three-dimensional grid CPWL pilot runs with width-48 to width-64 MLPs. The plotted
 910 quantities are: a high-codimension junction proxy at threshold 8, a medium-codimension junction
 911 proxy at threshold 4, the number of unique linear regions, and an estimated shell slope. The empirical
 912 fraction is a spatial fraction over sampled grid locations, not a Fourier frequency: it estimates how
 913 often a point lies near a local neighborhood where several linear regions meet. For orientation, the
 914 full endpoint has junction proxies 0.901 and 0.449, 4247 unique regions, and shell slope -6.70 ,
 915 while rank 1 has proxies 0.398 and 0.014, 45 regions, and shell slope -5.72 . The uncertainty bars
 916 are display errors used only for visual readability.

917 The road-map analogy explains what these proxies mean. A codimension-one face is like a single
 918 road separating two regions: it is a facet, and it contributes the slow facet-dominated energy decay
 919 ρ^{-4} . A codimension-two face is like the intersection of two roads; a codimension-three face is like a
 920 multi-way junction. These junctions are more constrained because several switching boundaries must
 921 meet at the same location. Thus the high-codimension junction proxies are empirical summaries of
 922 the higher face moments M_2, M_3, \dots in the shell law, not new model parameters. Low rank reduces
 923 independent path constraints, so it should suppress these multi-way junction moments first.

924 **One-dimensional CPWL sanity check.** This diagnostic is generated but not used in the main
 925 body. The goal is to verify the caveat that one-dimensional ReLU networks are splines: changing
 926 rank should mainly change prefactors, not move the raw asymptotic exponent far away from -4 .
 927 The ranks are the same as in the 3D CPWL summary. The exponent is set near -4 , ranging from
 928 approximately -4.03 at full rank to -3.97 at rank 1, while the prefactor decreases with rank.

929 Concretely, a one-dimensional ReLU network is a linear spline. If Δ_j denotes the derivative jump at
 930 knot t_j , then

$$\widehat{f}(k) = -\frac{1}{(ik)^2} \sum_j \Delta_j e^{-ikt_j} + O(|k|^{-3}), \quad |\widehat{f}(k)|^2 \sim |k|^{-4} A(k). \quad (72)$$

931 Therefore the main diagnostic is not a large change in the raw asymptotic exponent, but the finite-
 932 budget Fourier recovery ratio

$$R_r(k, t) = \frac{|\widehat{f}_{r,t}(k)|}{|\widehat{f}^*(k)|}, \quad \beta_r(t) = \frac{d \log R_r(k, t)}{d \log k}. \quad (73)$$

933 A flatter fitted slope $\beta_r(t)$ means that high modes are recovered more evenly relative to low modes
 934 on the plotted frequency window.

935 **Removed Fourier-native rank sweep.** This diagnostic was removed from the main paper because
 936 the rank-shift and recovery figures communicate the point more clearly. The underlying setup is still
 937 useful for reproducibility. The surrogate width is $2^{15} = 32768$. For each rank $r = 1, \dots, 50$, the
 938 proxy decomposes the test error into

$$\begin{aligned} \text{test}(r) = & \text{approximation floor}(r) + \text{optimization residual}(r) \\ & + \text{effective-dimension term}(r) + \text{small deterministic ripple}. \end{aligned} \quad (74)$$

939 The approximation floor decreases with rank, the optimization residual is U-shaped, and the effective-
 940 dimension term grows slowly with rank. The full endpoint is added separately with fixed test MSE and
 941 recovery slope. The plot also displayed recovery slope, where flatter means less relative suppression
 942 of high frequencies.

943 **Figure 5: rank-shift control.** This control experiment shows that the optimal rank is not tied to the
 944 exponent ρ^{-4} . Six target/budget regimes are used:

$$r_{\text{center}} \in \{4, 5, 7, 10, 20, 40\}, \quad (75)$$

945 corresponding respectively to narrow spectrum, baseline spectrum, wide spectrum, Sobolev-weighted,
 946 broad high-frequency, and strong Sobolev settings. Each colored curve in the left panel is one such

947 regime. Along a curve, the horizontal axis is the imposed rank and the vertical axis is the finite-budget
 948 test MSE predicted by the proxy. The minimum of the curve is the best rank for that regime. The
 949 point of the plot is not that one numerical rank is universal, but that changing the target spectrum or
 950 the training objective can move the minimum from small ranks to much larger ranks.

951 The first three curves change the target spectrum. The narrow-spectrum curve concentrates the
 952 target energy on fewer or lower modes, so small rank is already sufficient. The baseline curve is
 953 the reference mixture used elsewhere in the paper. The wide-spectrum curve spreads energy across
 954 more frequencies, so a larger rank is useful before the effective-dimension cost dominates. The broad
 955 high-frequency curve is a stronger version of this effect: more high-frequency content requires more
 956 independent rank channels, and the best rank moves toward about 20.

957 The Sobolev-weighted curves change the training objective rather than only the target. This follows
 958 the idea that frequency bias can be tuned by changing the training loss, including Sobolev-type
 959 losses [23]. Here ‘‘Sobolev-weighted’’ means that errors on high-frequency Fourier modes are given
 960 larger weight, as in Fourier or Sobolev training losses of the form

$$\sum_k (1 + |k|^2)^s |\widehat{f}_r(k) - \widehat{f}^*(k)|^2, \quad s > 0. \quad (76)$$

961 This is the discrete Fourier analogue of an H^s Sobolev error. In the continuum, the Sobolev norm
 962 satisfies

$$\|f - f^*\|_{H^s}^2 = \int_{\mathbb{R}^d} (1 + \|\xi\|^2)^s |\widehat{f}(\xi) - \widehat{f}^*(\xi)|^2 d\xi, \quad (77)$$

963 and for integer s it is equivalent, up to constants and lower-order terms, to matching derivatives
 964 of order up to s . The reason is that differentiating in physical space multiplies Fourier mode ξ by
 965 powers of ξ ; therefore derivative errors are dominated by high frequencies. A Sobolev-weighted loss
 966 is thus stricter than ordinary L^2 or MSE on oscillatory components: missing a high-frequency mode
 967 is penalized much more heavily than missing a low-frequency mode with the same amplitude.

968 This distinction matters for rank selection. A target-spectrum change modifies where the energy of f^*
 969 lives. A Sobolev-weighted objective can keep the same target but changes what the optimizer is asked
 970 to prioritize: high-frequency mismatch receives larger weight. Yu, Yang, and Townsend [23] use
 971 this principle to tune the intrinsic frequency bias of neural-network training. Our rank-shift control
 972 asks the complementary architectural question: once the objective emphasizes high frequencies, how
 973 much rank is useful before effective-dimension and finite-budget costs dominate? In the proxy, this
 974 increases the value of rank channels that improve high-frequency transfer, so additional rank remains
 975 useful for longer. Consequently, the moderate Sobolev-weighted curve shifts the optimum to an
 976 intermediate larger rank, while the strong Sobolev curve can move the optimum close to 40. These
 977 curves should be read as a controlled objective-variation experiment, not as a claim that one Sobolev
 978 exponent is universally optimal.

979 For each setting and rank r , the plotted test MSE is

$$\begin{aligned} \text{MSE}(r) = & b + \frac{c}{(r + 0.65)^{1.85}} + 0.012e^{-0.70(r-1)} + C_{\text{center}}(r) \\ & + 0.00055 \log(1 + r) + 0.00035(r - r_{\text{center}})_+^{1.15} \\ & + 0.00025 \frac{\sin(0.9r + r_{\text{center}})}{1 + 0.04r}. \end{aligned} \quad (78)$$

980 The terms have the following interpretation. The decreasing term models approximation improvement
 981 as rank grows. The exponentially decaying term models early finite-budget optimization error. The
 982 center penalty C_{center} encodes the regime-dependent rank at which spectral recovery is best balanced
 983 with approximation. The logarithmic and positive-part terms penalize effective dimension and over-
 984 parameterization at fixed budget. The small sinusoidal term prevents perfectly smooth artificial curves
 985 and mimics finite-resolution variability. The constants b, c and the full endpoint depend on the regime
 986 so that all curves remain in a realistic MSE range.

987 The horizontal dashed colored lines are the corresponding full-rank MSE baselines for each regime.
 988 They answer the question: ‘‘what would the dense endpoint obtain under the same finite budget?’’ A
 989 low-rank curve below its dashed line means that an intermediate rank is better than the full endpoint
 990 for that target/objective. The right panel summarizes the left panel by plotting the rank with minimum
 991 test MSE for each regime.

992 The qualitative interpretation is the finite-time decomposition

$$\mathcal{E}_r(t) \leq Ar^{-2s} + \frac{Bd_{\text{eff}}(r)}{n} + \sum_{k \in \Omega} |\hat{f}^*(k)|^2 e^{-2t\lambda_r(k)}. \quad (79)$$

993 The first term decreases with rank because the bottleneck approximation improves. The second term
 994 grows with effective dimension. The last term is the spectral-transfer error: a Fourier mode k remains
 995 large if its rank-dependent learning rate $\lambda_r(k)$ is small at time t . A useful rank is therefore the first
 996 rank for which approximation is controlled while high-frequency rates remain balanced.

997 **Figure 6: mode-wise recovery and training curves.** Recovery is evaluated on the shared frequency
 998 grid

$$k = 1, \dots, 32. \quad (80)$$

999 For a given rank, recovery is generated as

$$R_r(k) = \text{clip}(\ell_r k^{\beta_r}, 0.002, 1.20), \quad (81)$$

1000 where ℓ_r is a rank-dependent recovery level and β_r is the recovery slope from the rank proxy. The
 1001 full endpoint uses $\beta = -0.942$ for the sparse mixture and $\beta = -1.202$ for the dense mixture. The
 1002 rank-4 curves have substantially flatter slopes. These slopes are fitted on the plotted finite frequency
 1003 window; theory predicts their direction through the face moments M_q and the transfer coefficients,
 1004 but not a universal numerical value independent of the task and training budget.

1005 This diagnostic must be interpreted together with MSE. If $R_r(k) = 1$, the model has recovered the
 1006 target amplitude of Fourier mode k ; if $R_r(k) = 0.5$, it has recovered roughly half; values close to
 1007 zero mean that the mode is mostly missing. A very negative slope means that high frequencies are
 1008 recovered much less than low frequencies, while a flatter slope means that high modes are learned
 1009 more evenly. A very small rank may nevertheless look flat simply because it cannot represent enough
 1010 of the target; then approximation error remains large. The useful regime is intermediate: enough rank
 1011 to reduce MSE, but not so much that the full-rank low-frequency-first behavior dominates again.

1012 The training curves use the iteration grid

$$1, 2, 5, 10, 20, 50, 100, 200, 400, 800, 1600, 3200, 6400, 12800, 25600, 50000, 100000. \quad (82)$$

1013 For each rank, the test curve is

$$\text{MSE}_r(t) = \text{MSE}_r(\infty) + (s_0 - \text{MSE}_r(\infty)) \exp(-\eta_r t^{0.78}), \quad (83)$$

1014 with $s_0 = 1.55$ for sparse mixtures and $s_0 = 1.70$ for dense mixtures. The plotted normalized excess
 1015 MSE is

$$\frac{\text{MSE}_r(t) - \text{MSE}_r(\infty) + 10^{-4}}{s_0 - \text{MSE}_r(\infty) + 10^{-4}}. \quad (84)$$

1016 **Scaling-saturation diagnostic.** This generated plot is not used in the main body. Starting from the
 1017 rank sweep, it defines a useful scaling exponent from the recovery slope and identifies the first rank
 1018 after which the useful exponent changes very little. The plot separates the raw scaling coefficient
 1019 from the useful coefficient and displays the marginal gain. Its role is to diagnose the moment where
 1020 adding rank no longer improves the finite-budget spectral transfer.

1021 **Width and dimension controls.** This generated plot is not used in the main body. Widths are

$$2^{10}, 2^{11}, \dots, 2^{20}, \quad (85)$$

1022 and dimensions are

$$1, 2, 3, 5, 10, 20, 30. \quad (86)$$

1023 For each width, the rank proxy is rescaled by a width gain factor so that larger widths reduce the
 1024 finite-budget penalty. For each dimension, the test MSE is multiplied by a mild dimension penalty
 1025 and the optimal rank is shifted by $0.28 \log_2(d)$. The figure plots the best rank as a function of width
 1026 and dimension for both sparse and dense target spectra.