

# Weight Anisotropy in Mean-Field Theory: Learning on Isotropic Data

author names withheld

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

## Abstract

Neural networks efficiently learn isotropic data distributions with low-dimensional target structure where fixed kernel limits fail. We trace this advantage to input feature selection (IFS): networks develop strong weight anisotropy along task-relevant coordinates. While standard Mean-Field (MF) theory captures the onset of feature learning, it tracks only first moments and thus misses IFS and underestimates post-transition generalisation. We introduce MF-ARD, augmenting MF with a single additional set of order parameters for coordinate-wise precisions. MF-ARD successfully captures the sharp generalisation transitions of finite-width networks.

## 1. Introduction

While finite-width networks can discover low-dimensional targets in isotropic data via feature learning (FL) [7, 11, 12], fixed kernels [14, 17] suffer from the curse of dimensionality. Statistical physics and Mean-Field (MF) theory provide a rigorous framework for modeling this FL regime [3, 5, 9, 10, 13, 16, 19–22]. However, as outlined above, standard MF structurally misses **input feature selection (IFS)** because it lacks the second-moment tracking required to model coordinate-dependent regularisation. We resolve this by proposing **MF-ARD**, which integrates automatic relevance determination [18, 27] into the MF framework. By elevating coordinate-wise precisions to self-consistent order parameters, MF-ARD explicitly models the strong weight anisotropy observed in SGLD-trained networks. This minimal extension is sufficient to quantitatively predict the sharp generalisation transitions that standard MF underestimates. Extended related work is deferred to Section A.

## 2. Background: Mean field theory

### 2.1. SGLD setup

For analytic tractability, we use SGLD, the limit of full-batch GD with weight decay and injected Gaussian noise. ([26, 29], see Section C.1 for details). Given an NN described by  $f_{\theta}(\mathbf{x}_{\mu})$ , a dataset  $\mathcal{D} = \{(\mathbf{x}_{\mu}, \mathbf{y}_{\mu})\}_{\mu=1}^P$ , drawn from an input distribution  $q(\mathbf{x}, \mathbf{y})$ , the SGLD update equation is

$$\Delta\theta_{i,t} = -\eta \left[ \lambda_i \theta_{i,t} + \nabla_{\theta_i} \left( \frac{1}{P} \sum_{\mu=1}^P (f_{\theta}(\mathbf{x}_{\mu}) - \mathbf{y}_{\mu})^2 \right) \right] + \sqrt{2T\eta} \xi_{i,t}, \quad \xi_{i,t} \sim \mathcal{N}(0, 1), \quad (1)$$

with full-batch gradients on parameters  $\theta_i$ , injected Gaussian noise  $\xi$ , learning rate  $\eta$ , weight decay  $\lambda_i$ , and noise strength set by  $T = 2\kappa^2$ . Additionally, we focus on two-layer fully-connected

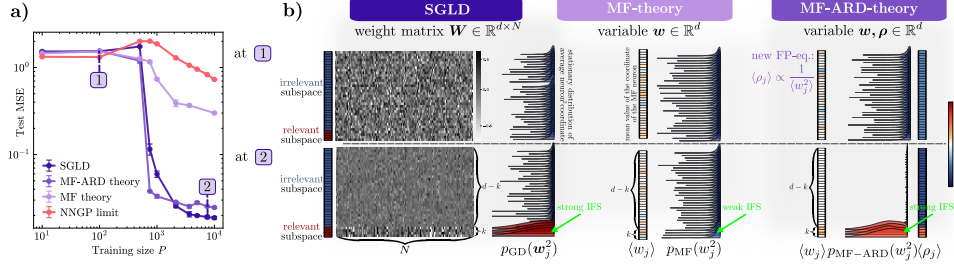


Figure 1: **a)** Test MSE vs. training-set size  $P$ . MF-ARD successfully captures the sharp generalization transition of SGLD, which standard MF theory underestimates and the NNGP limit misses entirely. Points **1** and **2** mark pre- and post-transition states. **b)** Stationary distributions of the squared weights,  $p(w_j^2)$ . At **2**, **SGLD** exhibits strong Input Feature Selection (IFS): the  $k$  relevant coordinates (red) develop large variances, while the  $d - k$  irrelevant ones (blue) remain suppressed. Standard **MF theory** tracks only coordinate means, yielding weak IFS. By introducing a fixed-point equation for coordinate variances, **MF-ARD** accurately reproduces the strong variance separation in  $p(w_j^2)$  observed in SGLD, explaining the post-transition performance gains.

networks with input dimension  $d$ , hidden width  $N$ , input weights  $\mathbf{w}_i \in \mathbb{R}^d$ , output weights  $a_i$ , output dimension 1, and nonlinearity  $\phi: f_\theta(\mathbf{x}) = \frac{1}{N^\gamma} \sum_{i=1}^N a_i \phi(\mathbf{w}_i^\top \mathbf{x})$ ,  $w_{ij} \sim \mathcal{N}(0, \frac{\sigma_w^2}{d})$ ,  $a_i \sim \mathcal{N}(0, \sigma_a^2)$ . The initial parameters are set with Gaussians. The scale factor  $N^{-\gamma}$  on the output layer enables interpolation between  $\gamma = 1/2$  (NTK scaling) and  $\gamma = 1$  (mean field scaling). Using the explicit form of  $f_\theta$  to rewrite the stationary distribution yields (see Section C.2):

### SGLD-posterior

$$\begin{aligned}
 -p_{\text{GD}}(\{\mathbf{w}_i\}_{i=1}^N, \mathbf{a} | \mathcal{D}) &= \frac{1}{2\sigma_a^2} \sum_{i=1}^N a_i^2 + \frac{d}{2\sigma_w^2} \sum_{i=1}^N \|\mathbf{w}_i\|^2 + \frac{1}{2\kappa^2 N^{2\gamma}} \sum_{i=1}^N a_i^2 \Sigma(\mathbf{w}_i) \\
 &+ \frac{1}{2\kappa^2 N^{2\gamma}} \sum_{i \neq i'} a_i a_{i'} G(\mathbf{w}_i, \mathbf{w}_{i'}) - \frac{1}{\kappa^2 N^\gamma} \sum_{i=1}^N a_i J_y(\mathbf{w}_i) + \text{const.},
 \end{aligned} \tag{2}$$

The posterior is shaped by three competing forces: the self-energy  $\Sigma(\mathbf{w}) = \frac{1}{P} \sum_\mu \phi(\mathbf{w}^\top \mathbf{x}_\mu)^2$  regularises neuron magnitude; the data coupling  $J_y(\mathbf{w}) = \frac{1}{P} \sum_\mu [\phi(\mathbf{w}^\top \mathbf{x}_\mu) y(\mathbf{x}_\mu)]$  rewards target alignment and breaks symmetry; and the interaction kernel  $G(\mathbf{w}_i, \mathbf{w}_{i'}) = \frac{1}{P} \sum_\mu \phi(\mathbf{w}_i^\top \mathbf{x}_\mu) \phi(\mathbf{w}_{i'}^\top \mathbf{x}_\mu)$  mediates neuron cooperation.

## 2.2. MF Theory: Removing off-diagonal couplings

To make the posterior in Equation (2) tractable, we use self-consistent MF theory (related to the MF theory in [22, 23]), which replaces the correlated posterior  $p(\mathbf{W} | \mathcal{D})$  with a factorised approximation  $p(\mathbf{W} | \mathcal{D}) \approx \prod_{i=1}^N p_{\text{MF}}(\mathbf{w}_i)$  where each neuron is independent. Each neuron interacts not with all others but with their average behavior, the mean field  $\langle f \rangle$ :  $\sum_{i' \neq i} a_{i'} G(\mathbf{w}_i, \mathbf{w}_{i'}) \stackrel{\text{MF}}{\approx} \frac{1}{P} \sum_\mu \phi(\mathbf{w}^\top \mathbf{x}_\mu) N^\gamma \langle f(\mathbf{x}) \rangle$ . We require  $\langle f(\mathbf{x}) \rangle = N^{1-\gamma} \cdot \mathbb{E}_{(\mathbf{w}, a) \sim p(\mathbf{w}, a)} [a \phi(\mathbf{w}^\top \mathbf{x})]$  for self-

consistency. Expanding  $\langle f(\mathbf{x}) \rangle = \sum_A m_A \chi_A(\mathbf{x})$  in an orthonormal basis  $\{\chi_A\}$  with feature coefficients  $m_A = \mathbb{E}_{\mathbf{x}}[\langle f(\mathbf{x}) \rangle \chi_A(\mathbf{x})]$  yields:

**MF theory (fixed point equations)**

$$\begin{aligned} -\ln p_{\text{MF}} &= \frac{a^2}{2\sigma_a^2} + \frac{d}{2\sigma_w^2} \sum_{j=1}^d w_j^2 + \frac{a^2}{2\kappa^2 N^{2\gamma}} \Sigma(\mathbf{w}) - \frac{a}{\kappa^2 N^\gamma} \left( J_Y(\mathbf{w}) - \sum_A m_A J_A(\mathbf{w}) \right) \\ m_A &= N^{1-\gamma} \langle a J_A(\mathbf{w}) \rangle_{\text{PMF}} \quad \forall A, \end{aligned} \quad (3)$$

where  $p_{\text{MF}} \rightarrow p(\mathbf{w}, a | \{m_A\}, \mathcal{D})$ ,  $\mathbf{w} \in \mathbb{R}^d$ ,  $a \in \mathbb{R}$  and  $J_A(\mathbf{w}) = \mathbb{E}_{\mathbf{x}}[\phi(\mathbf{w}^\top \mathbf{x}) \chi_A(\mathbf{x})]$ .

MF replaces the intractable  $N \times d$  weight matrix problem with a self-consistent problem for a single  $d$ -dimensional vector. These equations are solved by iteration, computing  $\Sigma(\mathbf{w})$ ,  $J_A(\mathbf{w})$ ,  $J_Y(\mathbf{w})$  directly from the finite training dataset of size  $P$  (Section G).

### 3. Theory I: Input Feature Selection as Implicit Regularisation

We focus on  $k$ -sparse parity,  $y(\mathbf{x}) = \chi_S(\mathbf{x})$  with  $S = \{1, \dots, k\}$ , as the canonical isotropic test case (Section A). Prior work [1, 2, 4, 24] established that weight-matrix anisotropy is central for generalisation; we give a simplified analytical derivation showing how it arises during SGLD training and shapes the stationary distribution. Let  $\mathbf{w}_j(t) \in \mathbb{R}^N$  denote the  $j$ -th row of  $\mathbf{W}(t) \in \mathbb{R}^{d \times N}$  and write  $\mathbf{w}_j^2(t) := \|\mathbf{w}_j(t)\|_2^2$ . Applying Itô's lemma to (1) yields the *exact* evolution

$$\frac{d}{dt} \langle \mathbf{w}_j^2 \rangle = -2\eta(\lambda \langle \mathbf{w}_j^2 \rangle + \langle \hat{\rho}_j \mathbf{w}_j^2 \rangle) + 2T\eta N. \quad (4)$$

$\hat{\rho}_j(t) := \frac{\langle \mathbf{w}_j, \nabla_{\mathbf{w}_j} \rangle}{\mathbf{w}_j^2}$  acts as a coordinate-wise effective decay: when  $\hat{\rho}_j < 0$ , the descent direction aligns with  $\mathbf{w}_j$ , partially cancelling  $\lambda$ . On sparse tasks,  $\hat{\rho}_j(t)$  separates sharply (Figure 3d): for  $j \in S^c$  the gradient carries no consistent component along  $\mathbf{w}_j$ , so  $\hat{\rho}_j \approx 0$  and the row experiences full decay; for  $j \in S$  the loss induces systematic alignment ( $\hat{\rho}_j < 0$ ) prior to the test-error drop, providing an anti-decay contribution that lets  $\mathbf{w}_j^2$  persist or grow. Setting  $\frac{d}{dt} \langle \mathbf{w}_j^2 \rangle = 0$  in (4) gives the exact balance  $\lambda \langle \mathbf{w}_j^2 \rangle_\infty + \langle \hat{\rho}_j \mathbf{w}_j^2 \rangle_\infty = TN$ , so a separation in  $\langle \hat{\rho}_j \rangle_\infty$  between  $S$  and  $S^c$  produces an anisotropic stationary distribution (Figure 3b, Figure 1a).

**Definition 1 (Input feature selection)** *A network undergoes IFS if, as epochs  $\rightarrow \infty$ , the weight-anisotropy ratio satisfies:  $R_{\mathbf{W}}^\infty = \frac{\sqrt{d-k} \langle \mathbf{w}_{j \in S}^2 \rangle_\infty}{\sqrt{k} \langle \mathbf{w}_{j \in S^c}^2 \rangle_\infty} \gg 1$ .*

Section D.1 confirms empirically that low test error tightly correlates with weight concentration in the relevant subspace (Figure 6).

### 4. Theory II: MF model of IFS

#### 4.1. FL in the simple MF model

MF theory interprets FL as *self-consistent symmetry breaking* in  $p_{\text{MF}}$ . Below  $(P_c, \kappa_c)$ , the only stable FP of Equation (3) is  $m_A = 0 \forall A$  and  $p_{\text{MF}}$  collapses to the Gaussian prior, regularised toward an isotropic high-entropy state by the prior and self-energy  $\Sigma(\mathbf{w})$ . Above  $(P_c, \kappa_c)$ , the neuron-data coupling  $J_Y$  acts as an external field that rewards alignment of  $\phi(\mathbf{w}^\top \mathbf{x})$  with  $y(\mathbf{x})$ , destabilising

$m_A = 0$  (see Theorem 11 for the explicit  $\kappa_c$ ). For  $y(\mathbf{x}) = \chi_S(\mathbf{x})$  the order parameter is  $m_S$ , and the FP equations (3) reduce to (Section F.5)  $m_S = \frac{N^{1-2\gamma}}{\kappa^2} (1 - m_S) \left\langle \frac{J_S(\mathbf{w})^2}{\sigma_a^{-2} + \Sigma(\mathbf{w}) / (\kappa^2 N^{2\gamma})} \right\rangle_{\mathbf{w} \sim p(\mathbf{w} | m_S)}$ . The transition occurs when  $\omega_0 > 1$  (Figure 4): SGLD and MF stay at  $m_S = 0$  until  $P_c$  then jump abruptly, while NNGP rises only smoothly.

#### 4.2. Problem: simple MF underestimates IFS

Figure 1a reveals the gap: post-transition, SGLD generalises far faster than MF. Standard MF parameterises  $p_{\text{MF}}(\mathbf{w})$  through *first-moment* order parameters  $\{m_A\}$  only, so the prior term  $\frac{d}{2\sigma_w^2} \sum_j w_j^2$  assigns equal precision to every  $j$  regardless of how  $\langle w_j^2 \rangle_{p_{\text{MF}}}$  varies between  $S$  and  $S^c$ . Tracking  $\{m_A\}$  predicts the *onset* of the transition (Figure 4) but not the post-transition gain, which Section 3 traces to symmetry breaking in the second moments  $\langle w_j^2 \rangle$ . We formalise this obstruction:

**Theorem 2** *Let  $y(\mathbf{x}) = \chi_S(\mathbf{x})$  and  $\{w_j\}_{j=1}^d$  a solution of (3) with  $\mathbf{w}_S = \{w_j\}_{j \in S}$ ,  $\mathbf{w}_{S^c} = \{w_j\}_{j \in S^c}$ . The plain-MF FP cannot develop strong weight anisotropy:  $R_{\mathbf{w}}^{\text{MF}} = \sqrt{d-k} \|\mathbf{w}_S\|^2 / (\sqrt{k} \|\mathbf{w}_{S^c}\|^2) = 1 + O(1/d)$ . See Section F.6 for the proof.*

#### 4.3. A minimal MF model of IFS: tracking coordinate variances

We extend MF to track the row-wise variances  $w_j^2$  as explicit order parameters. The goal is to match the SGLD stationary law (3): assuming the gradient-weight alignment behaves as a radial fixed force  $\langle \hat{\rho}_j \mathbf{w}_j^2 \rangle_\infty \approx \rho_j^r \langle \mathbf{w}_j^2 \rangle_\infty$  and using the MF independence assumption  $\langle \mathbf{w}_j^2 \rangle_\infty \approx N \langle w_j^2 \rangle_{\text{MF}}$ , (3) reduces to

$$\lambda \langle w_j^2 \rangle_{\text{MF}} + \rho_j^r \langle w_j^2 \rangle_{\text{MF}} = T, \quad (5)$$

i.e. the likelihood gradient acts as an effective potential  $\frac{\rho_j^r}{2} \|\mathbf{w}_j\|_2^2$  in coordinate  $j$  (linear response; see Figure 3c for empirical evidence). We therefore enlarge the state space minimally:  $\{(w_j, m_A)\}_{j=1}^d \rightarrow \{(w_j, m_A), (w_j^2, \rho_j)\}_{j=1}^d$ .

#### 4.4. Automatic Relevance Determination MF theory

We endow  $p_{\text{MF}}(\mathbf{w})$  with a coordinate-wise precision  $\rho_j$  via the conditional Gaussian  $p(w_j | \rho_j) = \mathcal{N}(0, \rho_j^{-1})$ , automatic relevance determination (ARD) in Bayesian NNs [18, 27], with a conjugate Gamma prior on  $\rho_j$ .

**Theorem 3** *Assume  $p(\mathbf{w}_j | \rho_j) = \mathcal{N}(0, \rho_j^{-1} \mathbb{I})$  and  $p(\rho_j) = \Gamma(\alpha_0, \beta_0)$ . Enforcing  $\rho_j \rightarrow d/\sigma_w^2$  when the data signal vanishes fixes  $\beta_0 = \alpha_0 \sigma_w^2 / d$ , giving the self-consistent FP  $\rho_j = (\alpha_0 + \frac{N}{2}) / (\frac{\alpha_0}{d} + \frac{N}{2} \langle w_j^2 \rangle_{p_{\text{ARD}}})$ . See Section C.7 for the proof.*

##### MF-ARD theory (fixed point equations)

$$-\ln p_{\text{ARD}} = \frac{a^2}{2\sigma_a^2} + \frac{1}{2} \sum_{j=1}^d \rho_j w_j^2 + \frac{a^2}{2\kappa^2 N^{2\gamma}} \Sigma(\mathbf{w}) - \frac{a}{\kappa^2 N^\gamma} \left( J_{\mathcal{Y}}(\mathbf{w}) - \sum_A m_A J_A(\mathbf{w}) \right) \quad (6)$$

$$m_A = N^{1-\gamma} \langle a J_A(\mathbf{w}) \rangle_{p_{\text{ARD}}} \quad \forall A \quad \rho_j = \frac{\alpha_0 + \frac{N}{2}}{\frac{\alpha_0}{d} + \frac{N}{2} \langle w_j^2 \rangle_{p_{\text{ARD}}}} \quad (7)$$

where  $p_{\text{ARD}} \rightarrow p_{\text{ARD}}(\mathbf{w}, a | \{m_A, \rho_i\}, \mathcal{D})$  and  $\mathbf{w} \in \mathbb{R}^d$ ,  $\boldsymbol{\rho} \in \mathbb{R}^d$ ,  $a \in \mathbb{R}$ .

The map  $\langle w_j^2 \rangle \mapsto \rho_j \mapsto p_{\text{ARD}}$  realises the IFS feedback: relevant coordinates acquire larger  $\langle w_j^2 \rangle$ , which lowers  $\rho_j$  via (7) and further reduces shrinkage; irrelevant coordinates do the opposite. Neurons remain independent, but the enlarged order-parameter set now captures SGLD’s coordinate-dependent variance structure.

#### 4.5. Connection of MF-ARD to IFS in SGLD

At stationarity, MF-ARD matches the SGLD distribution in both first moments (inherited from MF) and second moments up to  $O(\alpha_0/N)$ :

**Theorem 4** *Under the linear-response assumption  $\langle \hat{\rho}_j \mathbf{w}_j^2 \rangle_\infty \approx \rho_j^r \langle \mathbf{w}_j^2 \rangle_\infty$ , defining  $\rho_{\text{eff},j} := (\lambda + \rho_j^r)/T$ , and assuming  $\alpha_0 \ll N$ : SGLD:  $\rho_{\text{eff},j} = \frac{N}{\langle \mathbf{w}_j^2 \rangle_\infty}$ , MF-ARD:  $\rho_j = \frac{1}{\langle w_j^2 \rangle_{p_{\text{MF-ARD}}}} (1 + O(\alpha_0/N))$ . See Section F.7 for the proof.*

We thus identify  $\rho_j$  with the total effective SGLD precision (the factor  $N$  is absorbed by (5), and  $\alpha_0 \ll N$  throughout our experiments). Relevant features acquire strong negative alignment and large weights; irrelevant ones retain high precision, yielding an anisotropic stationary distribution from isotropic inputs, captured explicitly by MF-ARD’s static FP.

## 5. Results

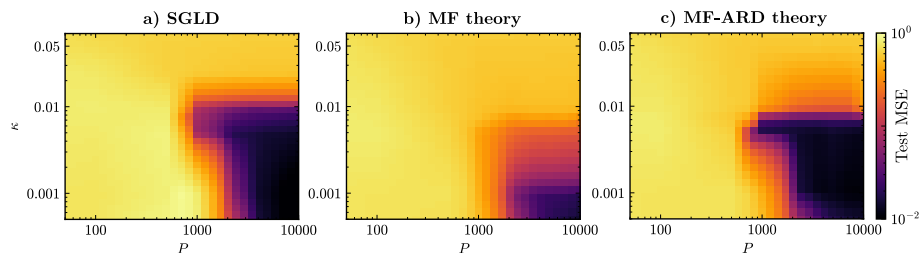


Figure 2: **Generalisation phase diagram (Test MSE) vs. dataset size  $P$  and noise  $\kappa$  for  $k$ -sparse parity target  $y(\mathbf{x}) = \chi_S(\mathbf{x})$  with  $S = \{0, 1, 2, 3\}$  in  $d = 35$ . a) SGLD, b) MF, c) MF-ARD: All panels show a transition from high to low test error as  $P$  increases or  $\kappa$  decreases. However, the plain MF theory strongly underestimates the sharpness of the transition while the MF-ARD theory largely reproduces the shape/location of this boundary. Training details are in Section H.**

Figure 2 shows test-MSE heatmaps over  $(P, \kappa)$  for SGLD, plain MF, and MF-ARD on  $k$ -sparse parity. Plain MF detects *when* learning starts but yields a diffuse, shifted boundary; MF-ARD, with the single added set  $\{\rho_j\}$ , tracks both the location and sharpness of the SGLD boundary and reproduces the “helpful noise” kink near  $\kappa = 0.05$ . The same picture holds for a single-index Hermite target [7] (Figure 8, Appendix).

**Limitations.** MF-ARD’s diagonal precision restricts feature selection to axis-aligned subspaces; rotated targets would require a low-rank covariance extension. The analysis uses the SGLD stationary distribution and is currently limited to two-layer networks; minibatch SGD and deeper architectures are left to future work.

## References

- [1] Emmanuel Abbe, Elisabetta Cornacchia, Jan Hazla, and Christopher Marquis. An Initial Alignment between Neural Network and Target is Needed for Gradient Descent to Learn. In *Proceedings of the 39th International Conference on Machine Learning*, pages 33–52. PMLR, June 2022. URL <https://proceedings.mlr.press/v162/abbe22a.html>. ISSN: 2640-3498.
- [2] Emmanuel Abbe, Enric Boix-Adsera, and Theodor Misiakiewicz. SGD learning on neural networks: leap complexity and saddle-to-saddle dynamics, August 2023. URL <http://arxiv.org/abs/2302.11055>. arXiv:2302.11055 [cs].
- [3] R Aiudi, R Pacelli, P Baglioni, A Vezzani, R Burioni, and P Rotondo. Local kernel renormalization as a mechanism for feature learning in overparametrized convolutional neural networks. *Nature Communications*, 16(1):568, 2025.
- [4] Jimmy Ba, Murat A. Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. High-dimensional Asymptotics of Feature Learning: How One Gradient Step Improves the Representation, May 2022. URL <http://arxiv.org/abs/2205.01445>. arXiv:2205.01445 [stat].
- [5] P Baglioni, R Pacelli, R Aiudi, F Di Renzo, A Vezzani, R Burioni, and P Rotondo. Predictive power of a bayesian effective action for fully connected one hidden layer neural networks in the proportional limit. *Physical Review Letters*, 133(2):027301, 2024.
- [6] Boaz Barak, Benjamin Edelman, Surbhi Goel, Sham Kakade, Eran Malach, and Cyril Zhang. Hidden progress in deep learning: Sgd learns parities near the computational limit. *Advances in Neural Information Processing Systems*, 35:21750–21764, 2022.
- [7] Alberto Bietti, Joan Bruna, Clayton Sanford, and Min Jae Song. Learning single-index models with shallow neural networks, 2022. URL <http://arxiv.org/abs/2210.15651>.
- [8] Blake Bordelon and Cengiz Pehlevan. Self-consistent dynamical field theory of kernel evolution in wide neural networks. *Advances in Neural Information Processing Systems*, 35:32240–32256, 2022.
- [9] Blake Bordelon and Cengiz Pehlevan. Dynamics of finite width kernel and prediction fluctuations in mean field neural networks. *Advances in Neural Information Processing Systems*, 36:9707–9750, 2023.
- [10] Michael Celentano, Chen Cheng, and Andrea Montanari. The high-dimensional asymptotics of first order methods with random data. *arXiv preprint arXiv:2112.07572*, June 2025. doi: 10.48550/arXiv.2112.07572. URL <https://arxiv.org/abs/2112.07572>.
- [11] Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural networks can learn representations with gradient descent. In *Proceedings of Thirty Fifth Conference on Learning Theory*, pages 5413–5452. PMLR, 2022. URL <https://proceedings.mlr.press/v178/damian22a.html>. ISSN: 2640-3498.

- [12] Amit Daniely and Eran Malach. Learning parities with neural networks. *Advances in Neural Information Processing Systems*, 33:20356–20365, 2020.
- [13] Jessica N Howard, Ro Jefferson, Anindita Maiti, and Zohar Ringel. Wilsonian renormalization of neural network gaussian processes. *Machine Learning: Science and Technology*, 6(2): 025038, 2025.
- [14] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- [15] Yiwen Kou, Zixiang Chen, Quanquan Gu, and Sham M. Kakade. Matching the Statistical Query Lower Bound for Sparse Parity Problems with Sign Stochastic Gradient Descent, December 2024. URL <http://arxiv.org/abs/2404.12376>. arXiv:2404.12376 [cs].
- [16] Clarissa Lauditi, Blake Bordelon, and Cengiz Pehlevan. Adaptive kernel predictors from feature-learning infinite limits of neural networks. *arXiv preprint arXiv:2502.07998*, 2025.
- [17] Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S. Schoenholz, Jeffrey Pennington, and Jascha Sohl-Dickstein. Deep neural networks as gaussian processes, 2018. URL <https://arxiv.org/abs/1711.00165>.
- [18] David JC MacKay. Bayesian non-linear modeling for the prediction competition. In *Maximum Entropy and Bayesian Methods: Santa Barbara, California, USA, 1993*, pages 221–234. Springer, 1996.
- [19] Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33): E7665–E7671, 2018.
- [20] Andrea Montanari and Pierfrancesco Urbani. Dynamical decoupling of generalization and overfitting in large two-layer networks. *arXiv preprint arXiv:2502.21269*, 2025.
- [21] Rosalba Pacelli, Sebastiano Ariosto, Mauro Pastore, Francesco Ginelli, Marco Gherardi, and Pietro Rotondo. A statistical mechanics framework for bayesian deep neural networks beyond the infinite-width limit. *Nature Machine Intelligence*, 5(12):1497–1507, 2023.
- [22] Zohar Ringel, Noa Rubin, Edo Mor, Moritz Helias, and Inbar Seroussi. Applications of statistical field theory in deep learning, 2025. URL <https://arxiv.org/abs/2502.18553>.
- [23] Noa Rubin, Inbar Seroussi, and Zohar Ringel. Grokking as a first order phase transition in two layer networks. *arXiv preprint arXiv:2310.03789*, May 2024. URL <https://arxiv.org/abs/2310.03789>. Preprint.
- [24] Zhenmei Shi, Junyi Wei, and Yingyu Liang. A theoretical analysis on feature learning in neural networks: Emergence from inputs and advantage over fixed features. *arXiv preprint arXiv:2206.01717*, 2022.

- [25] Taiji Suzuki, Denny Wu, Kazusato Oko, and Atsushi Nitanda. Feature learning via mean-field Langevin dynamics: classifying sparse parities and beyond. *Advances in Neural Information Processing Systems*, 36:34536–34556, December 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/hash/6cc321baf0a8611b1d1bdbd18822667b-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2023/hash/6cc321baf0a8611b1d1bdbd18822667b-Abstract-Conference.html).
- [26] Yee Whye Teh, Alexandre Thiéry, and Sebastian Vollmer. Consistency and fluctuations for stochastic gradient langevin dynamics, 2015. URL <https://arxiv.org/abs/1409.0578>.
- [27] Michael E Tipping. Sparse bayesian learning and the relevance vector machine. *Journal of machine learning research*, 1(Jun):211–244, 2001.
- [28] Alexander van Meegen and Haim Sompolinsky. Coding schemes in neural networks learning classification tasks. 16(1):3354, 2025. ISSN 2041-1723. doi: 10.1038/s41467-025-58276-6. URL <https://doi.org/10.1038/s41467-025-58276-6>.
- [29] Max Welling and Yee Whye Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th International Conference on Machine Learning (ICML 2011)*, ICML’11, pages 681–688, Bellevue, Washington, USA, June 2011. Omnipress. ISBN 978-1-4503-0619-5. doi: 10.5555/3104482.3104568. URL <https://dl.acm.org/doi/10.5555/3104482.3104568>.

## LLM Usage

LLMs were used to aid in checking and writing proofs, and also for reviewing the main text and to write code.

## Appendix A. Related work: Data-dependent feature learning and mean-field theories

We focus on  $k$ -sparse parity and sparse multi-index models because these isotropic tasks provide the sharpest stress test for feature learning (FL): the input distribution carries no anisotropy that would help a fixed kernel, so NTK/NNGP methods pay a sample-complexity exponent in  $d$  while finite-width networks do not [6, 11, 12, 15]. A rich line of gradient-based analyses proves that networks *do* recover the relevant subspace on such tasks [2, 4, 6, 11], but these results establish *that* learning happens rather than *predicting* generalisation error across the full  $(P, \kappa)$  plane. Our goal is a tractable theory whose fixed points predict the complete phase diagram (Figure 2), including the location and sharpness of the  $P_c$  transition. Several mean-field-style frameworks go beyond the kernel regime (Table A): classical MF and its Langevin extensions [19, 25]; DMFT-style theories tracking  $O(P^2T^2)$  kernel order parameters [8, 9, 16]; DMFT in the proportional limit for single-index models ( $k = 1$ , with  $k \geq 2$  deferred) [20]; and Bayesian/replica approaches [21, 28]. Despite their power, none has demonstrated a sample-complexity transition  $P_c$  on high-dimensional isotropic tasks (Table A). Two barriers explain this. First, DMFT tracks  $O(P^2T^2)$  order parameters requiring Monte Carlo solvers [8, 16], which becomes prohibitive at the dataset sizes ( $P \sim 10^3$ – $10^4$ ) needed to resolve the transition on sparse parity. Second, each existing closure is organised around the observable natural to its limit—sample-time kernels, latent-subspace projections, or replica  $Q$ -matrices—not around the input-coordinate row energies  $\{\|W_{j,:}\|_2^2\}_{j=1}^d$  that govern feature selection on isotropic inputs. MF-ARD addresses both gaps: it augments standard MF with only  $d$  coordinate-wise precisions  $\{\rho_j\}$ , the minimal closure that exposes IFS while remaining valid for any  $P, N$ . We use ARD [18, 27] as the mechanism that turns coordinate-wise second moments into self-consistent order parameters, not as a pruning heuristic.

Theory	#OP	Regime	Data tested	$P_c$ on isotropic
Rubin et al. [23]	$O(1)$	Bayesian, $N \rightarrow \infty$	cubic, mod. add.	No
Bordelon and Pehlevan [8]	$O(P^2T^2)$	$T, P = O_N(1)$	small CIFAR, Gaussian	No
Lauditi et al. [16]	$O(P^2T^2)$	$N \rightarrow \infty, \mu P$	CIFAR subsets, Gaussian	No
Montanari and Urbani [20]	$O(T^2)$	$n/d \rightarrow \alpha, m \rightarrow \infty$	$k=1$ single-index Gaussian	$k \geq 2$ deferred
Pacelli et al. [21]	$O(1)$	Bayesian, $P/N \rightarrow \alpha$	MNIST, CIFAR (proportional)	No
van Meegen and Sompolinsky [28]	$O(P)$	Bayesian, $N \rightarrow \infty$	MNIST, CIFAR	No
<b>MF-ARD (ours)</b>	<b><math>O(d)</math></b>	<b>any <math>P, N</math></b>	<b>high-<math>d</math> isotropic (<math>k</math>-parity, <math>k</math>-index)</b>	<b>Yes (Fig. 6)</b>

Table 1: Mean-field-style theories of feature learning. “#OP” = order-parameter count; “ $P_c$  on isotropic” = whether a sample-complexity transition on a high- $d$  isotropic task has been demonstrated.

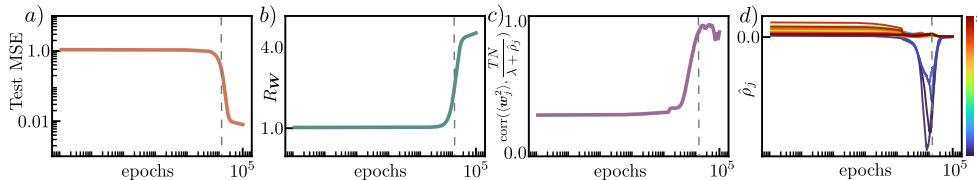


Figure 3: A 1-hidden layer ReLU network ( $N = 512$ ) is trained on  $k = 4$  sparse parity in  $d = 35$ . **a)** Test error versus epochs. **b)** The weight anisotropy ratio  $R_W$  increases sharply coincident with the drop in test error. **c)** Validation of the stationary energy balance (3), comparing measured  $w_j^2$  with the theoretical stationary prediction by plotting the mean correlation of  $w_j^2$  with  $\frac{TN}{\lambda + \rho_j}$ . **d)** Evolution of  $\hat{\rho}_j$ : support features ( $j \in S$ ) become strongly negative prior to the test error drop, while non-support features ( $j \in S^c$ ) remain near zero (colors indicate coordinate indices).

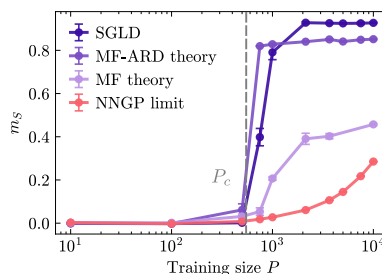


Figure 4: Feature coefficient  $m_S$  vs.  $P$  for a  $k$ -sparse parity target  $y(\mathbf{x}) = \chi_S(\mathbf{x}) = \prod_{j \in S} x_j$  with index  $S = \{0, 1, 2, 3\}$  in  $d = 35$  with a ReLU network ( $N = 512$ , see Section C.3; Section H for more details). SGLD, MF and MF-ARD exhibit a phase transition at  $P_c$ , NNGP grows smoothly, no phase-transition or FL.

## Appendix B. Main text figures

## Appendix C. Additional Background

### C.1. Equivalence of Langevin Dynamics and Bayesian Inference

In this section, we provide a detailed derivation for the equivalence between the stationary distribution of a network trained with Stochastic Gradient Langevin Dynamics (SGLD) and the posterior distribution of a corresponding Bayesian model.

**Network and Priors** We analyze a two-layer neural network with the functional form:

$$f(\mathbf{x}) = \frac{1}{N\gamma} \sum_{i=1}^N a_i \phi(\mathbf{w}_i^\top \mathbf{x}). \quad (8)$$

The parameters are drawn from independent Gaussian priors, which corresponds to choosing a specific prior in a Bayesian model. The initializations are given by:

- **Weights:**  $w_{ij} \sim \mathcal{N}(0, g_w^2)$  with  $g_w^2 = \frac{\sigma_w}{d}$ . This implies a prior probability  $p(\mathbf{w}_i) \propto \exp(-\frac{1}{2g_w^2} \|\mathbf{w}_i\|^2)$ .
- **Amplitudes:**  $a_i \sim \mathcal{N}(0, g_a^2)$  with  $g_a^2 = \sigma_a^2$ . This implies a prior probability  $p(a_i) \propto \exp(-\frac{1}{2g_a^2} \|a_i\|^2)$ .

Let  $\theta_i = (\mathbf{w}_i, a_i)$  denote the parameters for the  $i$ -th neuron.

**Stochastic Gradient Langevin Dynamics** We consider a full-batch Gradient Descent (GD) update for a parameter set  $\theta_i$ , which includes a weight decay term with coefficient  $\gamma$  and injected isotropic Gaussian noise  $\xi_t \sim \mathcal{N}(0, I)$ . This algorithm is known as Stochastic Gradient Langevin Dynamics (SGLD):

$$\Delta\theta_{i,t} := \theta_{i,t+1} - \theta_{i,t} \quad (9)$$

$$= -\eta(\gamma\theta_{i,t} + \nabla_{\theta_i} \ell(f_{\theta})) + \sqrt{2T\eta} \xi_t. \quad (10)$$

Here,  $\eta$  is the learning rate,  $T$  is a scalar temperature that controls the noise magnitude, and  $\ell(f_{\theta}) = \frac{1}{P} \sum_{\mu=1}^P (y_{\mu} - f(\mathbf{x}_{\mu}))^2$  is the mean squared error loss over the dataset of size  $P$ .

In the continuous-time limit ( $\eta \rightarrow 0$ ), this discrete update equation corresponds to a Langevin stochastic differential equation. The stationary distribution of this process, reached as  $t \rightarrow \infty$ , is given by the Gibbs-Boltzmann distribution:

$$p(\theta_i) \propto \exp\left(-\frac{1}{T} \left(\frac{\gamma}{2} \|\theta_i\|^2 + \ell(f_{\theta})\right)\right) \quad (11)$$

We can separate the weight decay terms  $\gamma$  into two separate parameters for weights and amplitudes,  $\gamma_w$  and  $\gamma_a$ , respectively. The stationary distribution for the parameters of a single neuron  $(\mathbf{w}_i, a_i)$  is then:

$$p(\mathbf{w}_i, a_i) \propto \exp\left(-\frac{1}{T} \left(\frac{\gamma_w}{2} \|\mathbf{w}_i\|^2 + \frac{\gamma_a}{2} \|a_i\|^2 + \frac{1}{P} \sum_{\mu=1}^P (y_{\mu} - f(\mathbf{x}_{\mu}))^2\right)\right). \quad (12)$$

**Bayesian Posterior Distribution** From a Bayesian perspective, we aim to find the posterior distribution of the parameters given the data  $\mathcal{D} = \{(\mathbf{x}_{\mu}, y_{\mu})\}_{\mu=1}^P$ . The posterior is given by Bayes' theorem:  $p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta)$ .

- The prior  $p(\theta)$  is defined by our choice of initialization:  $p(\theta) = \prod_i p(\mathbf{w}_i)p(a_i) \propto \exp\left(-\sum_i \left(\frac{1}{2g_w^2} \|\mathbf{w}_i\|^2 + \frac{1}{2g_a^2} \|a_i\|^2\right)\right)$ .
- The likelihood  $p(\mathcal{D}|\theta)$  is chosen to be a Gaussian distribution with variance  $\kappa^2$ , corresponding to the mean squared error loss:  $p(\mathcal{D}|\theta) \propto \exp\left(-\frac{1}{2\kappa^2 P} \sum_{\mu=1}^P (y_{\mu} - f(\mathbf{x}_{\mu}))^2\right)$ .

Combining these, the log-posterior is proportional to the negative of an energy function  $\mathcal{E}(\theta, \mathcal{D})$ . The posterior distribution for a single neuron's parameters is:

$$p(\mathbf{w}_i, a_i|\mathcal{D}) \propto \exp\left(-\left(\frac{1}{2g_w^2} \|\mathbf{w}_i\|^2 + \frac{1}{2g_a^2} \|a_i\|^2 + \frac{1}{2\kappa^2 P} \sum_{\mu=1}^P (y_{\mu} - f(\mathbf{x}_{\mu}))^2\right)\right). \quad (13)$$

**Identifying the Distributions** To establish the equivalence, we equate the functional forms of the SGLD stationary distribution and the Bayesian posterior. By comparing the exponents term-by-term, we find the following correspondences:

- **Weights:**  $\frac{\gamma_w}{2T} \|\mathbf{w}_i\|^2 = \frac{1}{2g_w^2} \|\mathbf{w}_i\|^2 \implies \frac{\gamma_w}{T} = \frac{1}{g_w^2}$
- **Amplitudes:**  $\frac{\gamma_a}{2T} \|a_i\|^2 = \frac{1}{2g_a^2} \|a_i\|^2 \implies \frac{\gamma_a}{T} = \frac{1}{g_a^2}$
- **Loss Term:**  $\frac{1}{TP} \sum_{\mu} \ell_{\mu} = \frac{1}{2\kappa^2 P} \sum_{\mu} \ell_{\mu} \implies T = 2\kappa^2$

This implies that the SGLD algorithm with temperature  $T$  effectively samples from a Bayesian posterior with data noise variance  $\kappa^2 = T/2$ , and with prior variances  $g_w^2 = T/\gamma_w$  and  $g_a^2 = T/\gamma_a$ .

**Final Update Equations** By substituting these relations back into the SGLD update rules, we obtain the dynamics for sampling from the desired posterior:

$$\Delta a_{t,i} = -\eta \left( \frac{T}{g_a^2} a_{t,i} + \nabla_{a_i} \left( \frac{1}{P} \sum_{\mu} \ell_{\mu} \right) \right) + \sqrt{2T\eta} \xi_{t,a} \quad (14)$$

$$\Delta \mathbf{w}_{t,i} = -\eta \left( \frac{T}{g_w^2} \mathbf{w}_{t,i} + \nabla_{\mathbf{w}_i} \left( \frac{1}{P} \sum_{\mu} \ell_{\mu} \right) \right) + \sqrt{2T\eta} \xi_{t,w}. \quad (15)$$

For training, this means we must set the temperature to  $T = 2\kappa^2$ .

## C.2. Derivation of (2)

We start from the negative log posterior and the two-layer model in Eq. (2.1):

$$-\ln p_{\text{GD}}(\boldsymbol{\theta} \mid \mathcal{D}) = \underbrace{\sum_l \frac{1}{2\sigma_l^2} \sum_j \|\boldsymbol{\theta}_{l,j}\|^2}_{-\ln p_{\text{prior}}} + \underbrace{\frac{1}{2\kappa^2 P} \sum_{\mu=1}^P (f_{\boldsymbol{\theta}}(\mathbf{x}_{\mu}) - y_{\mu})^2}_{-\ln p_{\text{L}}}.$$

For our two-layer network,  $f_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{N\gamma} \sum_{i=1}^N a_i \phi(\mathbf{w}_i^{\top} \mathbf{x})$ . Define the shorthand  $\phi_{i\mu} := \phi(\mathbf{w}_i^{\top} \mathbf{x}_{\mu})$ . Then the data term expands as

$$\begin{aligned} -\ln p_{\text{L}} &= \frac{1}{2\kappa^2 P} \sum_{\mu=1}^P \left( \frac{1}{N\gamma} \sum_{i=1}^N a_i \phi_{i\mu} - y_{\mu} \right)^2 \\ &= \frac{1}{2\kappa^2 P} \sum_{\mu=1}^P \left[ \frac{1}{N^2\gamma} \sum_{i=1}^N \sum_{j=1}^N a_i a_j \phi_{i\mu} \phi_{j\mu} - \frac{2}{N\gamma} \sum_{i=1}^N a_i \phi_{i\mu} y_{\mu} + y_{\mu}^2 \right]. \end{aligned} \quad (16)$$

Introduce the dataset-averaged quantities

$$\begin{aligned} \Sigma(\mathbf{w}) &:= \frac{1}{P} \sum_{\mu=1}^P \phi(\mathbf{w}^{\top} \mathbf{x}_{\mu})^2, & G(\mathbf{w}, \mathbf{w}') &:= \frac{1}{P} \sum_{\mu=1}^P \phi(\mathbf{w}^{\top} \mathbf{x}_{\mu}) \phi(\mathbf{w}'^{\top} \mathbf{x}_{\mu}), \\ J_{\mathbf{y}}(\mathbf{w}) &:= \frac{1}{P} \sum_{\mu=1}^P \phi(\mathbf{w}^{\top} \mathbf{x}_{\mu}) y_{\mu}. \end{aligned} \quad (17)$$

Using  $\sum_{i,j} = \sum_{i=j} + \sum_{i \neq j}$  in (16), we obtain

$$\begin{aligned}
 -\ln p_L = \frac{1}{2\kappa^2 N^{2\gamma}} & \left[ \sum_{i=1}^N a_i^2 \underbrace{\frac{1}{P} \sum_{\mu} \phi_{i\mu}^2}_{\Sigma(\mathbf{w}_i)} + \sum_{i \neq j} a_i a_j \underbrace{\frac{1}{P} \sum_{\mu} \phi_{i\mu} \phi_{j\mu}}_{G(\mathbf{w}_i, \mathbf{w}_j)} \right] \\
 & - \frac{1}{\kappa^2 N^\gamma} \sum_{i=1}^N a_i \underbrace{\frac{1}{P} \sum_{\mu} \phi_{i\mu} y_\mu}_{J_Y(\mathbf{w}_i)} + \frac{1}{2\kappa^2 P} \sum_{\mu=1}^P y_\mu^2. \tag{18}
 \end{aligned}$$

The prior part for our parameterization  $w_{ij} \sim \mathcal{N}(0, \sigma_w^2/d)$  and  $a_i \sim \mathcal{N}(0, \sigma_a^2)$  is

$$-\ln p_{\text{prior}} = \frac{1}{2\sigma_a^2} \sum_{i=1}^N a_i^2 + \frac{d}{2\sigma_w^2} \sum_{i=1}^N \|\mathbf{w}_i\|^2. \tag{19}$$

Combining (18) and (19), and discarding the  $\theta$ -independent constant, yields Eq. (2):

$$\begin{aligned}
 -\ln p_{\text{GD}}(\mathbf{W}, \mathbf{a} \mid \mathcal{D}) = \frac{1}{2\sigma_a^2} \sum_{i=1}^N a_i^2 + \frac{d}{2\sigma_w^2} \sum_{i=1}^N \|\mathbf{w}_i\|^2 + \frac{1}{2\kappa^2 N^{2\gamma}} \sum_{i=1}^N a_i^2 \Sigma(\mathbf{w}_i) \\
 + \frac{1}{2\kappa^2 N^{2\gamma}} \sum_{i \neq j} a_i a_j G(\mathbf{w}_i, \mathbf{w}_j) - \frac{1}{\kappa^2 N^\gamma} \sum_{i=1}^N a_i J_Y(\mathbf{w}_i) + \text{const.}
 \end{aligned}$$

This derivation holds for any nonlinearity  $\phi$ .

### C.3. $k$ -sparse parity target function

A  $k$ -sparse parity target function teacher is a single Walsh basis function on the Boolean hypercube. Let  $S \subseteq [d]$  with  $|S| = k$ . The Walsh function indexed by  $S$  is

$$\chi_S(\mathbf{x}) = \prod_{j \in S} x_j, \quad \mathbf{x} \in \{\pm 1\}^d,$$

(and, if  $\mathbf{x} \in [-1, 1]^d$ , one may use  $\chi_S(\mathbf{x}) = \prod_{j \in S} \text{sign}(x_j)$ ). The family  $\{\chi_S\}_{S \subseteq [d]}$  forms an orthonormal basis under the uniform product measure, i.e.  $\mathbb{E}[\chi_S(\mathbf{x})\chi_T(\mathbf{x})] = \mathbb{1}\{S = T\}$ . In our experiments the teacher is  $y(\mathbf{x}) = \chi_S(\mathbf{x})$ , when  $|S| = k$  this is exactly the  $k$ -parity problem.

### C.4. Relation to kernel eigenvalue outliers

The transition from  $m_S = 0$  to  $m_S > 0$  manifests as an outlier eigenvalue in the learned kernel. Recall that  $m_S = N^{1-\gamma} \langle a J_S(\mathbf{w}) \rangle$  measures the mean alignment between neurons and the target mode, where  $J_S(\mathbf{w}) = \mathbb{E}[\phi(\mathbf{w}^\top \mathbf{x}) \chi_S(\mathbf{x})]$  quantifies how well a single neuron with weights  $\mathbf{w}$  correlates with  $\chi_S$ . When  $m_S$  becomes non-zero, it indicates that neurons have collectively aligned

their weights to capture the target structure. Their  $J_S(\mathbf{w}_i)$  values have grown large. This alignment directly impacts the empirical kernel through

$$\widehat{R}_S^{(a)} = \frac{\chi_S^\top K \chi_S}{\text{tr}(K)} = \frac{\sum_{i=1}^N a_i^2 J_S(\mathbf{w}_i)^2}{\sum_{i=1}^N a_i^2 \Sigma(\mathbf{w}_i)}, \quad (20)$$

where  $K_{\mu\nu} = \frac{1}{N} \sum_{i=1}^N \phi(\mathbf{w}_i^\top \mathbf{x}_\mu) \phi(\mathbf{w}_i^\top \mathbf{x}_\nu)$  is the learned kernel. The numerator  $\sum_i a_i^2 J_S(\mathbf{w}_i)^2$  grows quadratically with the alignment strengths  $J_S(\mathbf{w}_i)$ , while the denominator (trace) remains roughly constant. In the kernel regime ( $m_S = 0$ ), neurons remain randomly oriented so  $J_S(\mathbf{w}_i) \sim O(N^{-1/2})$  and the ratio stays  $O(1/P)$ . However, when FL occurs ( $m_S > 0$ ), the enhanced  $J_S(\mathbf{w}_i)$  values cause  $\chi_S^\top K \chi_S$  to grow substantially, making  $\chi_S$  an outlying eigendirection of  $K$ . This anisotropic deformation, from the isotropic NNGP to a low-rank, plus isotropic structure, provides a direct spectral signature of the FL phase transition.

### C.5. The NNGP limit

The MF model can be simplified even further in the infinite-width limit with  $\gamma = 1/2$  (NTK-scaling). For this scaling the limit is well-behaved, the self-energy term in Equation (3) vanishes as well as any  $m_A$ -dependent tilt, resulting in  $p(\mathbf{w})$  collapsing to its prior (see Section F.1 for a *proof*).

#### NNGP limit (fixed point equations)

$$\begin{aligned} -\ln p_\infty &= \frac{d}{2\sigma_w^2} \|\mathbf{w}\|^2 \\ m_A^\infty &= \frac{\sigma_a^2}{\kappa^2} \left( \langle J_Y(\mathbf{w}) J_A(\mathbf{w}) \rangle_{p_\infty} \right. \\ &\quad \left. - \sum_B m_B^\infty \langle J_B(\mathbf{w}) J_A(\mathbf{w}) \rangle_{p_\infty} \right) \end{aligned} \quad (21)$$

where  $p_\infty \rightarrow p_\infty(\mathbf{w})$ ,  $\mathbf{w} \in \mathbb{R}^d$ .

**Kernel picture** We can define the kernel:  $K_{AB} := \mathbb{E}_{\mathbf{w} \sim p_\infty} [J_A(\mathbf{w}) J_B(\mathbf{w})]$  reducing the equations above to kernel ridge regression:

$$\begin{aligned} \mathbf{m}^\infty &= K (K + \tau \mathbb{1})^{-1} \mathbf{y}, \text{ where} \\ \tau &= \kappa^2 / \sigma_a^2, \quad y_A = \mathbb{E}_{\mathbf{x}} [y(\mathbf{x}) \chi_A(\mathbf{x})]. \end{aligned} \quad (22)$$

This NNGP limit represents the most restrictive limit in our approximation hierarchy. While MF theory eliminates inter-neuron interactions ( $w_{ij} \leftrightarrow w_{i'j}$ ), the NNGP limit additionally removes intra-neuron coordinate coupling ( $w_{ij} \leftrightarrow w_{ij'}$ ).

**No FL mechanism** In the infinite-width limit, there is *no* FL in the sense above. With no data-dependent term ( $J_Y$ ) to break symmetry, all neurons remain frozen at their prior. The feature coefficients  $m_A^\infty$  are now determined purely by the fixed kernel  $K_{AB} = \mathbb{E}_{\mathbf{w} \sim p_\infty} [J_A(\mathbf{w}) J_B(\mathbf{w})]$  (see Section C.6 how this connects to the usual kernel formulation of the NNGP limit).

### C.6. Connection to the standard NNGP formulation

The kernel representation in the function basis  $\{\chi_A\}$  can be transformed to recover the standard NNGP formulation in input space. We start with the function-basis kernel

$$K_{AB} = \mathbb{E}_{\mathbf{w} \sim p_\infty} [J_A(\mathbf{w})J_B(\mathbf{w})], \quad (23)$$

where  $J_A(\mathbf{w}) = \mathbb{E}_{\mathbf{x}}[\phi(\mathbf{w}^\top \mathbf{x})\chi_A(\mathbf{x})]$  projects the neuron's output onto basis function  $\chi_A$ . Expanding this definition:

$$K_{AB} = \mathbb{E}_{\mathbf{w}} \left[ \mathbb{E}_{\mathbf{x}}[\phi(\mathbf{w}^\top \mathbf{x})\chi_A(\mathbf{x})] \cdot \mathbb{E}_{\mathbf{x}'}[\phi(\mathbf{w}^\top \mathbf{x}')\chi_B(\mathbf{x}')] \right] \quad (24)$$

$$= \mathbb{E}_{\mathbf{w}} \mathbb{E}_{\mathbf{x}, \mathbf{x}'} \left[ \phi(\mathbf{w}^\top \mathbf{x})\phi(\mathbf{w}^\top \mathbf{x}')\chi_A(\mathbf{x})\chi_B(\mathbf{x}') \right] \quad (25)$$

$$= \mathbb{E}_{\mathbf{x}, \mathbf{x}'} \left[ \chi_A(\mathbf{x})\chi_B(\mathbf{x}') \cdot \underbrace{\sigma_a^2 \mathbb{E}_{\mathbf{w}}[\phi(\mathbf{w}^\top \mathbf{x})\phi(\mathbf{w}^\top \mathbf{x}')] }_{=: K(\mathbf{x}, \mathbf{x}')} \right], \quad (26)$$

where  $K(\mathbf{x}, \mathbf{x}')$  is the standard NNGP kernel in input space. The fixed-point equation  $m_A = \frac{\sigma_a^2}{\kappa^2}(\Xi_A - \sum_B K_{AB}m_B)$  with  $\Xi_A = \sum_S y_S K_{AS}$  becomes

$$m = K(K + \tau I)^{-1}y, \quad \tau = \kappa^2/\sigma_a^2. \quad (27)$$

To see this explicitly, consider data points  $\{\mathbf{x}_\mu\}_{\mu=1}^P$  with labels  $y_\mu$ . The predictor in the function basis is  $f(\mathbf{x}) = \sum_A m_A \chi_A(\mathbf{x})$ , while in the input basis it becomes  $f(\mathbf{x}_\mu) = \sum_{\nu=1}^P \alpha_\nu K(\mathbf{x}_\mu, \mathbf{x}_\nu)$  where  $\alpha = (K + \tau I)^{-1}y$ . The equivalence follows from the change of basis: if  $\chi_A(\mathbf{x}) = \delta_{\mathbf{x}, \mathbf{x}_A}$  (point evaluation basis), then  $K_{AB} = K(\mathbf{x}_A, \mathbf{x}_B)$  directly recovers the Gram matrix. For general orthogonal bases, the kernel ridge regression solution remains invariant under this transformation.

### C.7. Derivation of the ARD precision fixed point

Recall the ARD prior on per-coordinate precisions  $\boldsymbol{\rho} = (\rho_1, \dots, \rho_d)$  and weights:

$$p(\mathbf{w} \mid \boldsymbol{\rho}) = \prod_{j=1}^d \mathcal{N}(w_j \mid 0, \rho_j^{-1}), \quad p(\boldsymbol{\rho}) = \prod_{j=1}^d \Gamma(\rho_j \mid \alpha_0, \beta_0), \quad (28)$$

and the single-neuron MF-ARD action (Equation (6) in the main text)

$$-\ln p_{\text{ARD}}(\mathbf{w}, a \mid \{m_A\}, \boldsymbol{\rho}, \mathcal{D}) = \frac{a^2}{2\sigma_a^2} + \frac{1}{2} \sum_{j=1}^d \rho_j w_j^2 + \frac{a^2}{2\kappa^2 N^{2\gamma}} \Sigma(\mathbf{w}) - \frac{a}{\kappa^2 N^\gamma} \left( \mathcal{J}_y(\mathbf{w}) - \sum_A m_A J_A(\mathbf{w}) \right). \quad (29)$$

For a width- $N$  two-layer network, the joint action is the sum over neurons  $i = 1, \dots, N$  of Equation (29), and the (negative) log-evidence (free energy) for fixed  $\{m_A\}$  is

$$\mathcal{F}(\boldsymbol{\rho}; \{m_A\}) = -\ln Z(\boldsymbol{\rho}; \{m_A\}) - \ln p(\boldsymbol{\rho}) \quad \text{with} \quad Z(\boldsymbol{\rho}; \{m_A\}) = \int \prod_{i=1}^N d\mathbf{w}_i da_i e^{-\sum_{i=1}^N S_{\text{ARD}}(\mathbf{w}_i, a_i)}. \quad (30)$$

Stationarity of the negative log-evidence w.r.t.  $\rho_j$  gives the ARD FP:

$$0 = \frac{\partial \mathcal{F}}{\partial \rho_j} = \underbrace{\left\langle \frac{\partial}{\partial \rho_j} \sum_{i=1}^N S_{\text{ARD}}(\mathbf{w}_i, a_i) \right\rangle_{p_{\text{ARD}}}}_{\text{energy term}} - \underbrace{\frac{N}{2} \frac{1}{\rho_j}}_{\text{Gaussian normalizer}} - \underbrace{\frac{\partial \ln p(\boldsymbol{\rho})}{\partial \rho_j}}_{\text{prior term}} = \frac{1}{2} \sum_{i=1}^N \langle w_{ij}^2 \rangle_{p_{\text{ARD}}} - \frac{N}{2} \frac{1}{\rho_j} - \frac{\alpha_0 - 1}{\rho_j} + \beta_0, \quad (31)$$

where we used  $\partial S_{\text{ARD}}/\partial \rho_j = \frac{1}{2} \sum_i w_{ij}^2$  and  $\partial \ln \Gamma(\rho_j | \alpha_0, \beta_0)/\partial \rho_j = (\alpha_0 - 1)/\rho_j - \beta_0$ . By MF symmetry, all neurons are i.i.d. under  $p_{\text{ARD}}$ , so  $\sum_{i=1}^N \langle w_{ij}^2 \rangle_{p_{\text{ARD}}} = N \langle w_j^2 \rangle_{p_{\text{ARD}}}$ , and Equation (31) yields the closed-form FP update

$$\rho_j^* = \frac{\alpha_0 - 1 + \frac{N}{2}}{\beta_0 + \frac{N}{2} \langle w_j^2 \rangle_{p_{\text{ARD}}}} \stackrel{\text{MAP corr.}}{\approx} \frac{\alpha_0 + \frac{N}{2}}{\beta_0 + \frac{N}{2} \langle w_j^2 \rangle_{p_{\text{ARD}}}}, \quad (32)$$

where the MAP conjugacy correction (absorbing the  $-1$  in  $\alpha_0$ ) gives the form used in the main text (Equation (7)). With the scale-matching choice  $\beta_0 = \alpha_0/d$ , Equation (32) becomes

$$\rho_j^* = \frac{\alpha_0 + \frac{N}{2}}{\frac{\alpha_0}{d} + \frac{N}{2} \langle w_j^2 \rangle_{p_{\text{ARD}}}}. \quad (33)$$

### C.8. NNGP limit of the MF-ARD model

#### NNGP-limit of MF-ARD

$$\mathcal{S}_{\infty}^{\text{FL}}(\mathbf{w}|\rho) = \frac{1}{2} \sum_{j=1}^d \rho_j w_j^2, \quad p_{\infty, \rho}(\mathbf{w}) = \frac{1}{\mathcal{Z}} \exp(-\mathcal{S}_{\infty}^{\text{FL}}(\mathbf{w}|\rho)), \quad (34)$$

$$m_A^{\infty}(\rho) = \frac{\sigma_a^2}{\kappa^2} \left\langle (Jy(\mathbf{w}) - \sum_B m_B^{\infty}(\rho) J_B(\mathbf{w})) J_A(\mathbf{w}) \right\rangle_{p_{\infty, \rho}}, \quad \forall A. \quad (35)$$

$$0 = \frac{1}{2} \text{tr}[(A - (Ay)(Ay)^{\top}) \partial_{\rho_j} K_{\rho}] + \beta_0 - \frac{\alpha_0 - 1}{\rho_j}, \quad A = (K_{\rho} + \tau I)^{-1} \quad (36)$$

**Kernel picture** We can define the kernel :  $K_{\rho, AB} = \mathbb{E}_{\mathbf{w} \sim p_{\infty, \rho}}[J_A(\mathbf{w})J_B(\mathbf{w})]$  reducing the solution above to KRR

$$m^{\infty}(\rho) = K_{\rho}(K_{\rho} + \tau I)^{-1}y, \quad \tau = \kappa^2/\sigma_a^2 \quad (37)$$

**Usual (kernel) form** Let  $C_{\rho} = \text{diag}(\rho)^{-1}$  and  $K_{\rho}(\mathbf{x}, \mathbf{x}') = \sigma_a^2 \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, C_{\rho})}[\phi(\mathbf{w}^{\top} \mathbf{x})\phi(\mathbf{w}^{\top} \mathbf{x}')]$ . Then

$$m^{\infty}(\rho) = K_{\rho}(K_{\rho} + \tau I)^{-1}y, \quad K_{\rho}\chi_A = \lambda_A(\rho)\chi_A \Rightarrow m_A^{\infty}(\rho) = \frac{\lambda_A(\rho)}{\lambda_A(\rho) + \tau} y_A. \quad (38)$$

This shows that  $\rho$  rescales coordinates before the nonlinearity, so  $K_{\rho} = \mathbb{E}_{\mathbf{z}}[J_A(C_{\rho}^{1/2} \mathbf{z}) J_B(C_{\rho}^{1/2} \mathbf{z})]$  is a nonlinear deformation of the isotropic kernel. Any nontrivial change in  $\rho$  therefore produces an  $O(1)$  change in the Gaussian measure over  $\mathbf{w}$ , hence an  $O(1)$  change in  $K$  (it does not vanish with width). ARD improves spectral alignment by increasing  $\lambda_A(\rho)$  along task-aligned directions. When  $\lambda_A(\rho)$  crosses the scale  $\tau$ ,  $m_A^{\infty}$  exhibits a jump, yielding the leading-order FL effect at infinite width.

### C.9. Universality of IFS

Our IFS mechanism is largely agnostic to the specific choice of nonlinearity, provided the nonlinearity allows for feature learning. As shown in Figure 5 networks trained with Sigmoid activations exhibit the same coordinate-wise heavy-tailed distributions on relevant input dimensions as ReLU networks. While the population-level organization (sparse vs. redundant neurons) may differ as predicted by [28], the fundamental mechanism of discovering the relevant input subspace via IFS remains the same.

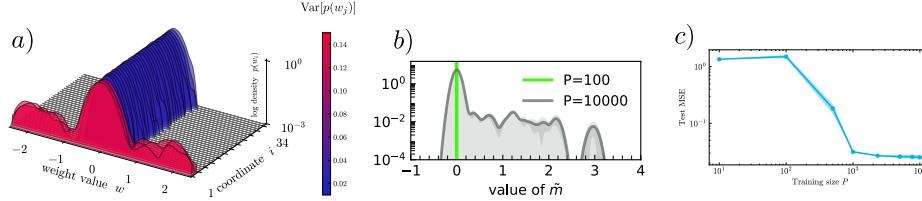


Figure 5: Similar to the ReLU results in the main text, SGLD training with Sigmoid activations leads to strong coordinate-wise symmetry breaking (high variance on relevant features  $j \in S$ ). This demonstrates that IFS is a fundamental mechanism for subspace discovery, distinct from the neuron-level coding schemes discussed in [28].

## Appendix D. Additional figures

### D.1. Empirical evidence for IFS

**An ensemble of trained NNs** In Figure 6, we draw random dataset sizes  $P \in [10, 20000]$  and noise levels  $\kappa_0 \in [5 \cdot 10^{-5}, 5 \cdot 10^{-2}]$ . We train the full NN for  $10^5$  epochs and color code the final test error. This confirms empirically that a minimal level of anisotropy is needed to generalise.

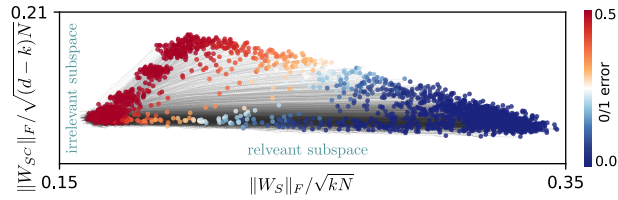


Figure 6: Each point corresponds to a one-hidden-layer  $N = 512$  network trained with SGLD on  $k$ -sparse parity, with training-set size  $P \in [10, 20000]$  and noise level  $\kappa_0 \in [5 \cdot 10^{-5}, 5 \cdot 10^{-2}]$ . The x-axis shows the normalized Frobenius norm of the first-layer weights restricted to the task-relevant coordinates  $S$ ,  $\|W_S\|_F / \sqrt{kN}$ , and the y-axis shows the corresponding norm on the irrelevant coordinates  $S^c$ ,  $\|W_{S^c}\|_F / \sqrt{(d-k)N}$ . Color denotes final test error. Low test error is achieved only when weight mass concentrates in the relevant subspace, illustrating that strong IFS  $R_W^\infty \gg 1$  correlates with generalisation in this isotropic setting.

**Causal validation** To show that IFS is the causal driver of generalisation, we introduce a counterfactual regularisation term to suppress IFS. We define the isotropy penalty:  $\ell_{\text{iso}}(\mathbf{W}) = \lambda_h \sum_{i=1}^N \text{Var}_{j \in [d]}(|w_{ji}|)$ . This term penalizes the variance of weight magnitudes across input dimensions  $j$  within each neuron  $i$ . Minimizing  $\ell_{\text{iso}}$  forces the network to maintain isotropic usage of inputs (i.e.,  $|w_{1i}| \approx |w_{2i}| \approx \dots$ ), effectively “switching off” the IFS mechanism while allowing other forms of coordination between neurons (including those discussed in [28]). Figure 7 clearly shows that growing regularisation causes the NN to generalise progressively worse. This confirms that IFS is not just a correlate of learning, but the necessary mechanism for generalisation in this regime.

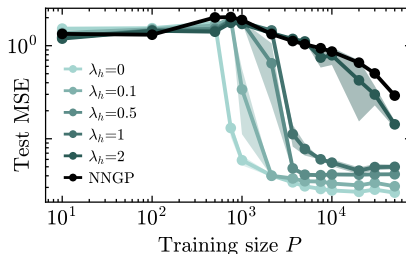


Figure 7: Test loss for  $k$ -sparse parity as a function of the regularisation  $\lambda_h$  that forces coordinates within a neuron to be homogeneous (preventing IFS). Increasing  $\lambda_h$  reduces performance till the NN performance degrades back to the Kernel (NNGP) limit. This confirms that coordinate-wise anisotropy is the necessary condition for sample-efficient learning in this setting.

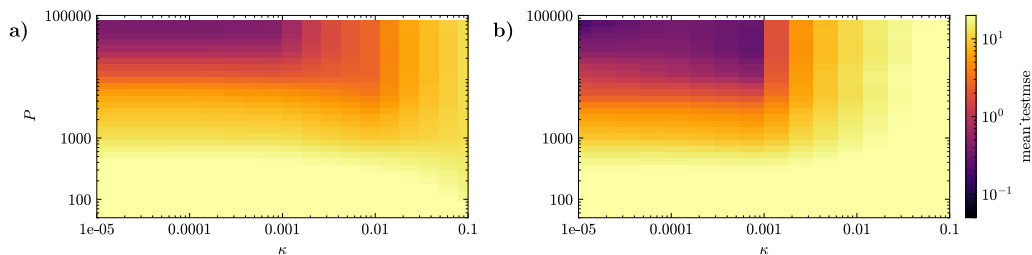


Figure 8: Same plot as Figure 2 a) SGLD, b) MF-ARD, but for a single-index model with Gaussian inputs  $\mathbf{x} \sim \mathcal{N}(0, \mathbb{1})$  and teacher  $y = \text{He}_p(\mathbf{w}^\top \mathbf{x})$ , where  $\mathbf{w}_j = 1/\sqrt{k}$  on a  $k$ -sized support and 0 otherwise, here  $d = 18$ ,  $k = 2$ ,  $p = 4$ .

## Appendix E. Comparison with recent mean field works

In this section, we clarify the relationship and fundamental differences between our MF-ARD framework and the recent work by van Meegen and Sompolinsky [28]. While both works utilize mean field approximations to study the posterior of neural networks, they address fundamentally different phenomena occurring at different levels of the network hierarchy.

### E.1. Coordinate-wise vs. Neuron-wise Sparsification

The most critical distinction lies in the domain where symmetry breaking occurs:

- **(Neuron-wise Sparsification):** [28] investigates how the population of  $N$  neurons organizes to represent the target. They find that the posterior distribution over the weight vectors  $\{\mathbf{w}_1, \dots, \mathbf{w}_N\}$  breaks permutation symmetry, resulting in a solution where only a subset of neurons ( $O(1)$  out of  $N$ ) strongly correlate with the target, while others remain agnostic or redundant. This is a *inter-neuron* phenomenon.
- **Our Work (Coordinate-wise Sparsification / IFS):** Our theory models the posterior of a *single representative neuron*  $\mathbf{w} = [w_1, \dots, w_d] \in \mathbb{R}^d$ . We investigate how the symmetry between input coordinates  $j = 1, \dots, d$  is broken. We find that for tasks with isotropic data but low-dimensional targets, the posterior marginals  $p(w_j)$  become heavy-tailed and high-variance only for task-relevant coordinates ( $j \in S$ ), while remaining Gaussian and narrow for irrelevant coordinates ( $j \notin S$ ). This is an *intra-neuron* phenomenon.

Standard mean field theories (including the base model in [28]) typically assume that the weight distribution is isotropic over input coordinates unless the prior enforces otherwise. Consequently, while [28] successfully predicts that some neurons become "active" and others "dormant," it does not explicitly model the self-reinforcing mechanism that allows an *active* neuron to selectively amplify specific input coordinates to overcome the curse of dimensionality.

### E.2. The Necessity of Coordinate Anisotropy: A Control Experiment

To rigorously prove that coordinate-wise symmetry breaking (the core of our MF-ARD theory) is the causal mechanism for beating the kernel regime, we performed a control experiment enforcing coordinate homogeneity.

We trained networks on the  $k$ -sparse parity task using a modified loss function that penalizes anisotropy within neurons, without penalizing sparsity across neurons:

$$\ell_{\text{total}} = \ell_{\text{MSE}} + \lambda_h \sum_{j=1}^N \left[ \frac{1}{d} \sum_{i=1}^d \left( |w_{ij}| - \left( \frac{1}{d} \sum_{k=1}^d |w_{kj}| \right) \right)^2 \right]. \quad (39)$$

**Result:** As  $\lambda_h$  increases, the learning curve of the neural network reverts to the NNGP learning curve (see Figure 7). Even if the network is allowed to sparsify at the neuron level (some neurons large, some small), suppressing the variance *between coordinates*  $w_{ij}$  prevents the network from learning the target efficiently.

This confirms that the performance gap between Finite-Width NNs and Kernels on isotropic data is driven by the mechanism modeled in this paper (IFS/Coordinate Anisotropy), a mechanism that is distinct from and complementary to the population-level phenomena described in [28].

## Appendix F. Proofs

### F.1. Kernel limit

**Theorem 5** *With  $\gamma = 1/2$  the infinite width limit has the following FP equations:*

$$m_A^\infty = \frac{\sigma_a^2}{\kappa^2} \left( \mathbb{E}_{\mathbf{w} \sim p_\infty} [J_A(\mathbf{w}) J_Y(\mathbf{w})] - \sum_B \mathbb{E}_{\mathbf{w} \sim p_\infty} [J_A(\mathbf{w}) J_B(\mathbf{w})] m_B^\infty \right). \quad (40)$$

**Proof** The proof proceeds by taking the  $N \rightarrow \infty$  limit of the self-consistency equations derived from the cavity method's free energy.

**1. Self-Consistency from Free Energy.** We begin with the MF free energy functional derived from the cavity method:

$$\mathcal{F}(\{m_A\}) = \frac{1}{2\kappa^2} \sum_A (y_A - m_A)^2 - N \ln \mathcal{Z}_1(\{m_A\}) \quad (41)$$

where  $\mathcal{Z}_1(\{m_A\}) = \int d\mathbf{w} e^{-\bar{S}_{\text{eff}}^{(\infty)}(\mathbf{w}; \{m_A\})}$  is the single-neuron partition function. The effective action for a single weight vector  $\mathbf{w}$ , after integrating out the amplitude  $a$ , is:

$$\bar{S}_{\text{eff}}^{(\infty)}(\mathbf{w}; \{m_A\}) = \frac{d}{2\sigma_w^2} \|\mathbf{w}\|^2 + \frac{1}{2} \ln \left( \frac{1}{\sigma_a^2} + \frac{\Sigma(\mathbf{w})}{N^{2\gamma}\kappa^2} \right) - \frac{(J_Y(\mathbf{w}) - \sum_A m_A J_A(\mathbf{w}))^2}{2\kappa^4 N^{2\gamma} \left( \frac{1}{\sigma_a^2} + \frac{\Sigma(\mathbf{w})}{N^{2\gamma}\kappa^2} \right)} \quad (42)$$

The equilibrium state is found by the stationarity condition  $\partial\mathcal{F}/\partial m_A = 0$ , which yields the self-consistency equation:

$$m_A = N^{1-\gamma} \langle \mu(\mathbf{w}) J_A(\mathbf{w}) \rangle_{p(\mathbf{w}|\{m_B\})} \quad (43)$$

where  $p(\mathbf{w}|\{m_B\}) = \frac{1}{\mathcal{Z}_1} e^{-\bar{S}_{\text{eff}}^{(\infty)}}$  is the posterior distribution on a single neuron's weights, and  $\mu(\mathbf{w})$  is the posterior mean of its amplitude  $a$ :

$$\mu(\mathbf{w}) = \frac{\frac{1}{\kappa^2 N^\gamma} (J_Y(\mathbf{w}) - \sum_B m_B J_B(\mathbf{w}))}{\frac{1}{\sigma_a^2} + \frac{\Sigma(\mathbf{w})}{N^{2\gamma}\kappa^2}} \quad (44)$$

**2. The Infinite-Width Limit with Critical Scaling.** To obtain a non-trivial limit as  $N \rightarrow \infty$ , we must set the scaling to the critical value  $\gamma = 1/2$ . For  $\gamma > 1/2$ , the prefactor  $N^{1-\gamma} \rightarrow 0$ , decoupling the neurons and leading to  $m_A \rightarrow 0$ . For  $\gamma < 1/2$ , the interaction term diverges. With  $\gamma = 1/2$ , the terms of order  $1/N$  in the effective action  $\bar{S}_{\text{eff}}^{(\infty)}$  vanish. Specifically:

$$\frac{\Sigma(\mathbf{w})}{N^{2\gamma}\kappa^2} = \frac{\Sigma(\mathbf{w})}{N\kappa^2} \xrightarrow{N \rightarrow \infty} 0 \quad (45)$$

As a result, the data-dependent and  $m_A$ -dependent terms in the exponent of  $p(\mathbf{w}|\{m_A\})$  vanish. The distribution over weights collapses to its prior:

$$p(\mathbf{w}|\{m_A\}) \xrightarrow{N \rightarrow \infty} p_\infty(\mathbf{w}) \propto \exp \left( -\frac{d}{2\sigma_w^2} \|\mathbf{w}\|^2 \right) \quad (46)$$

Simultaneously, the denominator in  $\mu(\mathbf{w})$  simplifies to  $1/\sigma_a^2$ . The expression for the posterior mean amplitude becomes:

$$\mu(\mathbf{w}) \xrightarrow{N \rightarrow \infty} \frac{\sigma_a^2}{\kappa^2 N^{1/2}} \left( J_{\mathcal{Y}}(\mathbf{w}) - \sum_B m_B J_B(\mathbf{w}) \right) \quad (47)$$

**3. Deriving the Linear System.** We now substitute these limiting forms back into the self-consistency equation, using  $m_A^\infty$  to denote the solution in this limit:

$$m_A^\infty = N^{1-1/2} \left\langle \left[ \frac{\sigma_a^2}{\kappa^2 N^{1/2}} \left( J_{\mathcal{Y}}(\mathbf{w}) - \sum_B m_B^\infty J_B(\mathbf{w}) \right) \right] J_A(\mathbf{w}) \right\rangle_{p_\infty(\mathbf{w})} \quad (48)$$

$$m_A^\infty = \frac{\sigma_a^2}{\kappa^2} \left\langle \left( J_{\mathcal{Y}}(\mathbf{w}) - \sum_B m_B^\infty J_B(\mathbf{w}) \right) J_A(\mathbf{w}) \right\rangle_{p_\infty(\mathbf{w})} \quad (49)$$

$$m_A^\infty = \frac{\sigma_a^2}{\kappa^2} \left( \mathbb{E}_{\mathbf{w} \sim p_\infty} [J_A(\mathbf{w}) J_{\mathcal{Y}}(\mathbf{w})] - \sum_B \mathbb{E}_{\mathbf{w} \sim p_\infty} [J_A(\mathbf{w}) J_B(\mathbf{w})] m_B^\infty \right) \quad (50)$$

Let us define the NNGP kernel matrix  $K$  with elements  $K_{AB} := \mathbb{E}_{\mathbf{w} \sim p_\infty} [J_A(\mathbf{w}) J_B(\mathbf{w})]$  and the data-kernel coupling vector  $\Xi$  with elements  $\Xi_A := \mathbb{E}_{\mathbf{w} \sim p_\infty} [J_A(\mathbf{w}) J_{\mathcal{Y}}(\mathbf{w})] = \sum_S y_S K_{AS}$ . The equation becomes a linear system for the vector  $m^\infty$ :

$$m^\infty = \frac{\sigma_a^2}{\kappa^2} (\Xi - K m^\infty) = \frac{\sigma_a^2}{\kappa^2} (K y - K m^\infty) \quad (51)$$

**4. Solution as Kernel Ridge Regression.** Rearranging the linear system, we get:

$$\left( I + \frac{\sigma_a^2}{\kappa^2} K \right) m^\infty = \frac{\sigma_a^2}{\kappa^2} K y \quad (52)$$

$$\left( \frac{\kappa^2}{2\sigma_a^2} I + \frac{1}{2} K \right) m^\infty = \frac{1}{2} K y \quad (53)$$

Let the ridge be  $\tau = \kappa^2/(2\sigma_a^2)$ . The equation is  $(\tau I + K)m^\infty = K y$ . Solving for  $m^\infty$  gives the final KRR solution:

$$m^\infty = K(K + \tau I)^{-1} y \quad (54)$$

This completes the proof. The solution depends only on the NNGP kernel  $K$ , which is fixed by the network architecture and priors, not the training data labels. This demonstrates the absence of FL in this specific infinite-width limit.  $\blacksquare$

## F.2. Integrating out $a$

As the action is quadratic in  $a$  we can integrate it out.

**Theorem 6** *The distribution is given by*

$$S_{\text{MF}}(\mathbf{w}, \{m_A\}) = \text{const.} + \frac{1}{2} \ln \left( \frac{1}{\sigma_a^2} + \frac{\Sigma(\mathbf{w})}{N^{2\gamma} \kappa^2} \right) - \frac{(J_{\mathcal{Y}}(\mathbf{w}) - \sum_A m_A J_A(\mathbf{w}))^2}{2N^{2\gamma} \kappa^4 \left( \frac{1}{\sigma_a^2} + \frac{\Sigma(\mathbf{w})}{N^{2\gamma} \kappa^2} \right)} \quad (55)$$

$$+ \frac{d}{2\sigma_w^2} \|\mathbf{w}\|^2 \quad (56)$$

$$p(a|\mathbf{w}) = \mathcal{N}(\mu(\mathbf{w}), \sigma^2(\mathbf{w})) \quad (57)$$

$$\sigma(\mathbf{w})^2 = \frac{1}{2\alpha} = \left( \frac{1}{\sigma_a^2} + \frac{\Sigma(\mathbf{w})}{N^{2\gamma}\kappa^2} \right)^{-1} \quad (58)$$

$$\mu(\mathbf{w}) = \sigma^2\beta = \frac{\beta}{2\alpha} = \frac{\frac{1}{N\gamma\kappa^2} (J\mathcal{Y}(\mathbf{w}) - \sum_A m_A J_A(\mathbf{w}))}{\left( \frac{1}{\sigma_a^2} + \frac{\Sigma(\mathbf{w})}{N^{2\gamma}\kappa^2} \right)} \quad (59)$$

**Proof** We start with the standard Gaussian integral

$$\int da d\mathbf{w} e^{-[\alpha a^2 - \beta \cdot a + c]} = \int d\mathbf{w} \sqrt{\frac{\pi}{\alpha}} e^{\frac{\beta^2}{4\alpha} - c} \quad (60)$$

We have

$$\mathcal{S}_{\text{MF}}(\mathbf{w}, a, \{m_A\}) = \frac{1}{2\sigma_a^2} \sum_{i=1}^N a_i^2 + \frac{d}{2\sigma_w^2} \sum_{i=1}^d \|\mathbf{w}_i\|^2 + \frac{a^2}{2\kappa^2 N^{2\gamma}} \Sigma(\mathbf{w}) \quad (61)$$

$$- \frac{a}{\kappa^2 N^\gamma} \left( J\mathcal{Y}(\mathbf{w}) - \sum_A m_A J_A(\mathbf{w}) \right) \quad (62)$$

For 1 neuron we get

$$\mathcal{S}_{\text{MF}}(\mathbf{w}, a, \{m_A\}) = a^2 \left( \frac{1}{2\sigma_a^2} + \frac{\Sigma(\mathbf{w})}{2\kappa^2 N^{2\gamma}} \right) - \frac{a}{\kappa^2 N^\gamma} \left( J\mathcal{Y}(\mathbf{w}) - \sum_A m_A J_A(\mathbf{w}) \right) + \frac{d}{2\sigma_w^2} \|\mathbf{w}\|^2 \quad (63)$$

with

$$\alpha := \frac{1}{2\sigma_a^2} + \frac{\Sigma(\mathbf{w})}{2N^{2\gamma}\kappa^2}, \quad (64)$$

$$\beta := \frac{J\mathcal{Y}(\mathbf{w}) - \sum_A m_A J_A(\mathbf{w})}{N^\gamma \kappa^2}. \quad (65)$$

$$c = \frac{d}{2\sigma_w^2} \|\mathbf{w}\|^2 \quad (66)$$

and  $\alpha(a - \frac{\beta}{2\alpha})^2 - \frac{\beta^2}{4\alpha} + c$  comparing to  $\frac{-(a-\mu)^2}{2\sigma^2}$  we get  $\alpha = \frac{1}{2\sigma^2}$  and  $\mu = \frac{\beta}{2\alpha}$  for the Gaussian identification. This gives

$$\int d\mathbf{w} \sqrt{\frac{\pi}{\alpha}} e^{\frac{\beta^2}{4\alpha} - c} = \int d\mathbf{w} e^{-[\frac{1}{2} \ln(\pi) - \frac{1}{2} \ln(\alpha) + \frac{\beta^2}{4\alpha} - c]} \quad (67)$$

where the exponent is identified as the effective action. ■

### F.3. MF- FP equation

**Lemma 7** Consider a single target mode  $y(\mathbf{x}) = \chi_S(\mathbf{x})$  and assume other overlaps vanish. Using the self-consistency  $m_A = N^{1-\gamma} \langle a J_A(\mathbf{w}) \rangle$  and the conditional mean  $\mu(\mathbf{w})$  above, we obtain

$$m_S = N^{1-\gamma} \langle \mu(\mathbf{w}) J_S(\mathbf{w}) \rangle = \frac{N^{1-2\gamma}}{\kappa^2} (1 - m_S) \left\langle \frac{J_S(\mathbf{w})^2}{\sigma_a^{-2} + \frac{\Sigma(\mathbf{w})}{\kappa^2 N^{2\gamma}}} \right\rangle_{\mathbf{w} \sim p(\mathbf{w}|m_S)}. \quad (68)$$

and in the  $\kappa \rightarrow 0$  limit it is

$$m_S \approx (1 - m_S) N \left\langle \frac{J_S(\mathbf{w})^2}{\Sigma(\mathbf{w})} \right\rangle_{\mathbf{w} \sim p(\mathbf{w}|m_S)} \quad (69)$$

**Proof** Consider a single target mode  $y(\mathbf{x}) = \chi_S(\mathbf{x})$  and assume other overlaps vanish. Using the self-consistency  $m_A = N^{1-\gamma} \langle a J_A(\mathbf{w}) \rangle$  and the conditional mean  $\mu(\mathbf{w})$  above, we obtain

$$m_S = N^{1-\gamma} \langle \mu(\mathbf{w}) J_S(\mathbf{w}) \rangle = \frac{N^{1-2\gamma}}{\kappa^2} (1 - m_S) \left\langle \frac{J_S(\mathbf{w})^2}{\sigma_a^{-2} + \frac{\Sigma(\mathbf{w})}{\kappa^2 N^{2\gamma}}} \right\rangle_{\mathbf{w} \sim p(\mathbf{w}|m_S)}.$$

In the small-noise regime  $\kappa \rightarrow 0$  (with  $\Sigma(\mathbf{w}) > 0$ ), the denominator simplifies as  $\sigma_a^{-2} + \frac{\Sigma(\mathbf{w})}{\kappa^2 N^{2\gamma}} \approx \frac{\Sigma(\mathbf{w})}{\kappa^2 N^{2\gamma}}$ , so that

$$m_S \approx (1 - m_S) N \left\langle \frac{J_S(\mathbf{w})^2}{\Sigma(\mathbf{w})} \right\rangle_{\mathbf{w} \sim p(\mathbf{w}|m_S)} \implies m_S \approx \frac{N \langle J_S(\mathbf{w})^2 / \Sigma(\mathbf{w}) \rangle}{1 + N \langle J_S(\mathbf{w})^2 / \Sigma(\mathbf{w}) \rangle}$$

By Cauchy-Schwarz,  $J_S(\mathbf{w})^2 \leq \Sigma(\mathbf{w})$ , so  $0 \leq m_S < 1$  unless all mass concentrates on perfectly aligned  $\mathbf{w}$ .  $\blacksquare$

### F.4. Solution to the fixed point equation

**Lemma 8** Consider  $\phi = \text{ReLU}$ . Let the inputs  $x_j \in \{\pm 1\}$  be i.i.d. and  $y(\mathbf{x}) = \chi_S(\mathbf{x}) = \prod_{j \in S} x_j$  with  $S = \{0, 1, \dots, k-1\}$ . We define  $R_k(\mathbf{w}) = \frac{J_S(\mathbf{w})^2}{\Sigma(\mathbf{w})}$ . It holds that

$$\mathbf{w}^* = \max_{\mathbf{w}} R_k(\mathbf{w}) \quad (70)$$

is given by  $\mathbf{w}^* = (\overbrace{\alpha, \dots, \alpha}^k, 0, \dots, 0)$ .

**Proof** Decompose  $w_S = \alpha \frac{\mathbb{1}_S}{\sqrt{k}} + u$ ,  $u \perp \mathbb{1}_S$  and keep arbitrary  $\mathbf{w}^C$ . Because the data distribution is invariant to permutations of the  $k$  coordinates in  $S$  the only  $S$ -dependent statistic that survives in  $\chi_S$ -weighted expectations is the sum  $s(\mathbf{x}) = \sum_{j \in S} x_j$ . Any component orthogonal to  $\mathbb{1}_S$  averages out in  $J_S$  but increases  $\Sigma(\mathbf{w})$  (by convexity of  $x \mapsto \phi(x)^2$ ). Likewise, weights in  $S^C$  contribute variance to  $\Sigma$  but contribute nothing to  $J_S$  (they are independent of  $\chi_S$  and average to zero under the sign symmetry). Thus the maximizer of  $R_k(\mathbf{w})$  lives in the span of  $\mathbb{1}_S$  and  $\mathbf{w}^C$ .  $\blacksquare$

**Lemma 9** Consider  $\phi = \text{ReLU}$ . Let the inputs  $x_j \in \{\pm 1\}$  be i.i.d. and  $y(\mathbf{x}) = \chi_S(\mathbf{x}) = \prod_{j \in S} x_j$  with  $|S| = k$ . We assume  $\mathbf{w}^*$  from Theorem 8. Then, the ratio of the squared neuron-target coupling  $J_S(\mathbf{w})^2$  to the neuron's self-energy  $\Sigma(\mathbf{w})$  is a constant  $R_k$  independent of the scale  $\alpha$ .

**Proof**

Let  $r \sim \text{Binomial}(k, 1/2)$  be the number of components  $x_j = -1$  for  $j \in S$ . The inner product is  $s(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} = \alpha \sum_{j \in S} x_j = \alpha(k - 2r)$ , and the target function is  $\chi_S(\mathbf{x}) = (-1)^r$ .

The neuron-target coupling  $J_S(\mathbf{w})$  is the expectation  $\mathbb{E}[\phi(\mathbf{w}^\top \mathbf{x})\chi_S(\mathbf{x})]$ .

$$J_S(\mathbf{w}) = \mathbb{E}[\alpha \cdot [k - 2r]_+ \cdot (-1)^r] \quad (71)$$

$$= \alpha \sum_{r=0}^k P(r) \cdot [k - 2r]_+ \cdot (-1)^r \quad (72)$$

$$= \alpha \cdot 2^{-k} \sum_{r=0}^{\lfloor (k-1)/2 \rfloor} \binom{k}{r} (k - 2r) (-1)^r \quad (73)$$

where  $[z]_+ = \max(0, z)$ , and the sum is restricted to the terms where  $k - 2r > 0$ .

The neuron's self-energy  $\Sigma(\mathbf{w})$  is the expectation  $\mathbb{E}[\phi(\mathbf{w}^\top \mathbf{x})^2]$ .

$$\Sigma(\mathbf{w}) = \mathbb{E}[(\alpha \cdot [k - 2r]_+)^2] \quad (74)$$

$$= \alpha^2 \sum_{r=0}^k P(r) \cdot [k - 2r]_+^2 \quad (75)$$

$$= \alpha^2 \cdot 2^{-k} \sum_{r=0}^{\lfloor (k-1)/2 \rfloor} \binom{k}{r} (k - 2r)^2 \quad (76)$$

We define the scale-independent constants  $D_k$  and  $C_k$ :

$$D_k := 2^{-k} \sum_{r=0}^{\lfloor (k-1)/2 \rfloor} \binom{k}{r} (k - 2r) (-1)^r \quad (77)$$

$$C_k := 2^{-k} \sum_{r=0}^{\lfloor (k-1)/2 \rfloor} \binom{k}{r} (k - 2r)^2 \quad (78)$$

such that  $J_S(\mathbf{w}) = \alpha D_k$  and  $\Sigma(\mathbf{w}) = \alpha^2 C_k$ . The ratio is then

$$R_k := \frac{J_S(\mathbf{w})^2}{\Sigma(\mathbf{w})} = \frac{(\alpha D_k)^2}{\alpha^2 C_k} = \frac{D_k^2}{C_k}$$

which is independent of  $\alpha$ , thus proving the proposition. ■

### E.5. Exact solution of the FP equation

To evaluate this expectation analytically, we employ a saddle-point (Laplace) approximation, a standard technique in the statistical mechanics of learning. Specifically, we assume that in the relevant noise regime, the posterior distribution of the weights concentrates sharply around its optimal direction  $\mathbf{w}^*$  (the mode). While a rigorous proof of posterior contraction in this specific high-dimensional, non-convex setting remains a challenging open problem for future work, this ansatz allows us to extract the leading-order macroscopic behavior. Under this saddle-point approximation, the expectation  $\langle \cdot \rangle_{\mathbf{w}|m_S}$  reduces to an evaluation at the optimal weight scale  $\alpha^*$ .

**Theorem 10** *Using the setup from Theorem 9, especially the proposed  $\mathbf{w}^*$ , the FP is given by*

$$m_S = (1 - m_S)NR_k \left( 1 - \frac{\sigma_a^{-2}}{A^*(m_S)} \right), \quad A^*(m_S) = \frac{-C_k + \sqrt{C_k^2 + 4 \left( \frac{dk}{\sigma_w^2} \right) N^{2\gamma} (1 - m_S)^2 D_k^2 \sigma_a^{-2}}}{2 \left( \frac{dk}{\sigma_w^2} \right) \kappa^2 N^{2\gamma}} \quad (79)$$

**Proof** From Theorem 7, we know:

$$m_S = \frac{N^{1-2\gamma}}{\kappa^2} (1 - m_S) \left\langle \frac{J_S(\mathbf{w})^2}{\sigma_a^{-2} + \Sigma(\mathbf{w})/(\kappa^2 N^{2\gamma})} \right\rangle_{\mathbf{w}|m_S}$$

Assuming the posterior distribution of weights is sharply peaked around the optimal direction given by the ansatz in Theorem 9, the expectation  $\langle \cdot \rangle_{\mathbf{w}|m_S}$  reduces to an evaluation at the optimal weight scale  $\alpha^*$ . The value of  $\alpha^*$  is determined by minimizing the single-neuron effective action  $\mathcal{S}_{\text{eff}}(\alpha)$ :

$$\mathcal{S}_{\text{eff}}(\alpha) = \underbrace{\frac{dk}{2\sigma_w^2} \alpha^2}_{\text{Weight Prior}} + \underbrace{\frac{1}{2} \ln A(\alpha)}_{\text{Normalizer}} - \underbrace{\frac{1}{2} \frac{[(1 - m_S)J_S(\alpha)]^2 / (\kappa^2 N^{2\gamma})}{A(\alpha)}}_{\text{Data Gain}}$$

where  $A(\alpha) = \sigma_a^{-2} + \Sigma(\alpha)/(\kappa^2 N^{2\gamma}) = \sigma_a^{-2} + \alpha^2 C_k / (\kappa^2 N^{2\gamma})$ .

Setting  $\frac{\partial \mathcal{S}_{\text{eff}}}{\partial \alpha} = 0$  yields a quadratic equation for the optimal value  $A^* = A(\alpha^*)$ :

$$\left( \frac{dk}{\sigma_w^2} \right) \kappa^2 N^{2\gamma} (A^*)^2 + C_k A^* - \frac{(1 - m_S)^2 D_k^2 \sigma_a^{-2}}{\kappa^2} = 0$$

The positive root of this equation gives the solution for  $A^*(m_S)$ :

$$A^*(m_S) = \frac{-C_k + \sqrt{C_k^2 + 4 \left( \frac{dk}{\sigma_w^2} \right) N^{2\gamma} (1 - m_S)^2 D_k^2 \sigma_a^{-2}}}{2 \left( \frac{dk}{\sigma_w^2} \right) \kappa^2 N^{2\gamma}}$$

We now substitute this back into the self-consistency equation. Using the definitions of  $J_S$ ,  $\Sigma$ , and  $A^*$ , we have  $(\alpha^*)^2 = \frac{\kappa^2 N^{2\gamma}}{C_k} (A^* - \sigma_a^{-2})$ .

$$m_S = (1 - m_S) \frac{N^{1-2\gamma}}{\kappa^2} \frac{J_S(\alpha^*)^2}{A^*} \quad (80)$$

$$= (1 - m_S) \frac{N^{1-2\gamma}}{\kappa^2} \frac{(\alpha^*)^2 D_k^2}{A^*} \quad (81)$$

$$= (1 - m_S) \frac{N^{1-2\gamma}}{\kappa^2} \frac{1}{A^*} \left[ \frac{\kappa^2 N^{2\gamma}}{C_k} (A^* - \sigma_a^{-2}) \right] D_k^2 \quad (82)$$

$$= (1 - m_S) N \frac{D_k^2}{C_k} \frac{A^* - \sigma_a^{-2}}{A^*} \quad (83)$$

$$= (1 - m_S) N R_k \left( 1 - \frac{\sigma_a^{-2}}{A^*(m_S)} \right) \quad (84)$$

This gives the final fixed-point equation for the order parameter  $m_S$ . ■

**Theorem 11** *Consider the fixed-point equation in the infinite  $P$ -limit*

$$m_S = (1 - m_S) N R_k \left( 1 - \frac{\sigma_a^{-2}}{A^*(m_S)} \right), \quad (85)$$

with  $A^*(m_S)$  given implicitly as the positive root of

$$\left( \frac{dk}{d\sigma_w^2} \right) \kappa^2 N^{2\gamma} (A^*)^2 + C_k A^* - (1 - m_S)^2 D_k^2 \sigma_a^{-2} = 0, \quad (86)$$

and define  $C := \frac{dk}{d\sigma_w^2}$ . Then the critical noise level

$$\kappa_c^2 = \frac{\sqrt{C_k^2 + 4C N^{2\gamma} D_k^2 \sigma_a^{-2}} - C_k}{2C \sigma_a^{-2} N^{2\gamma}} \quad (87)$$

marks a phase transition:

- (i) If  $\kappa^2 \geq \kappa_c^2$ , then  $A^*(0) \leq \sigma_a^{-2}$  and  $m_S = 0$  is a fixed point; in particular for  $\kappa^2 = \kappa_c^2$  we have  $A^*(0) = \sigma_a^{-2}$  and the only solution is  $m_S = 0$ .
- (ii) If  $\kappa^2 < \kappa_c^2$ , then  $A^*(0) > \sigma_a^{-2}$  and there exists a unique nontrivial solution  $m_S \in (0, 1)$ , thus the system exhibits symmetry breaking  $m_S : 0 \rightarrow m_S > 0$  as  $\kappa$  crosses  $\kappa_c$  from above.

**Proof** Set  $m_S = 0$  in (86) to get

$$C \kappa^2 N^{2\gamma} (A^*)^2 + C_k A^* - D_k^2 \sigma_a^{-2} = 0. \quad (88)$$

At the onset of FL the trivial fixed point  $m_S = 0$  changes stability precisely when the right-hand side of the FP map ceases to vanish at  $m_S = 0$ , i.e. when

$$N R_k \left( 1 - \frac{\sigma_a^{-2}}{A^*(0)} \right) = 0 \iff A^*(0) = \sigma_a^{-2}. \quad (89)$$

Plugging  $A^* = \sigma_a^{-2}$  into the quadratic and solving for  $\kappa^2$  gives

$$C \kappa_c^2 N^{2\gamma} \sigma_a^{-4} + C_k \sigma_a^{-2} - D_k^2 \sigma_a^{-2} = 0, \quad (90)$$

which, after rearrangement, yields

$$\kappa_c^2 = \frac{\sqrt{C_k^2 + 4 C N^{2\gamma} D_k^2 \sigma_a^{-2}} - C_k}{2 C \sigma_a^{-2} N^{2\gamma}}, \quad (91)$$

the stated expression.

For  $\kappa^2 > \kappa_c^2$  the quadratic gives  $A^*(0) \leq \sigma_a^{-2}$ , hence  $1 - \sigma_a^{-2}/A^*(0) \leq 0$  and the FP map evaluates to 0 at  $m_S = 0$ , so  $m_S = 0$  is a (and in fact the only) solution. For  $\kappa^2 < \kappa_c^2$  we have  $A^*(0) > \sigma_a^{-2}$  so the FP map at  $m_S = 0$  is strictly positive, continuity and the fact that the map is strictly decreasing in  $m_S$  (because  $(1 - m_S)$  and  $A^*(m_S)$  both decrease with  $m_S$ ) imply a unique intersection with the diagonal in  $(0, 1)$ , i.e. a unique  $m_S \in (0, 1)$  solves the FP equation. ■

We use the following notation. Let  $S \subset [d]$  with  $|S| = k$  denote the (unknown) support. For a weight coordinate  $w_j$  at a given outer iterate, write

$$v_j := \langle w_j^2 \rangle_p, \quad v'_j(\rho) := \mathbb{E}_{p_\rho}[w_j^2] \quad (92)$$

for the current and next second moments, respectively. For a vector of ARD precisions  $\rho \in \mathbb{R}_+^d$  define the explicit ARD map

$$\rho_j(v_j) = \frac{\alpha_0 + \frac{N}{2}}{\frac{\alpha_0}{d} + \frac{N}{2} v_j} =: \frac{A}{B(v_j)}, \quad A := \alpha_0 + \frac{N}{2}, \quad B(v) := \frac{\alpha_0}{d} + \frac{N}{2} v. \quad (93)$$

We write  $w_{-j} := (w_1, \dots, w_{j-1}, w_{j+1}, \dots, w_d)$  for all coordinates except  $j$  and write  $p_\rho$  for the posterior  $p_{\text{ARD}}$  to make the  $\rho$  dependence explicit.

**Assumption** Here, we will state the assumption that is needed to prove the theorem:  $\varepsilon$  **symmetry breaking towards  $S$** . There exists an outer iterate  $t_0 = \mathcal{O}(1)$  and a constant  $\varepsilon_0 > 0$  (independent of  $d$ ) such that

$$\min_{j \in S} v_j^{t_0} - \max_{j \notin S} v_j^{t_0} \geq c \varepsilon_0, \quad v_j^t := \langle w_j^2 \rangle_p \text{ at outer time } t. \quad (94)$$

We need to establish the following global bound.

**Lemma 12** Fix  $j \in [d]$  and  $w_{-j}$ . Assume  $\|x\|_\infty \leq 1$  (e.g.  $x \in \{\pm 1\}^d$ ) and  $\phi(z) = \max(0, z)$ , as well as  $m_S \leq 1$  Let

$$g(\mathbf{w}) = \frac{1}{2} \ln \left( \sigma_a^{-2} + \frac{\Sigma(\mathbf{w})}{\kappa^2 N^{2\gamma}} \right) - \frac{(J_y(\mathbf{w}) - m_S J_S(\mathbf{w}))^2}{2 \kappa^4 N^{2\gamma} \left( \sigma_a^{-2} + \frac{\Sigma(\mathbf{w})}{\kappa^2 N^{2\gamma}} \right)}. \quad (95)$$

Then there exists  $L_\star > 0$ , independent of  $d$  and  $w_{-j}$ , such that the map  $t \mapsto g(w_{-j}, t)$  satisfies

$$|\partial_t^2 g(w_{-j}, t)| \leq L_\star \quad \text{for a.e. } t \in \mathbb{R}. \quad (96)$$

Consequently, for all  $t \in \mathbb{R}$ ,

$$g(w_{-j}, t) \geq g(w_{-j}, 0) - \frac{L_\star}{2} t^2. \quad (97)$$

**Proof** Fix  $j \in [d]$  and  $w_{-j}$ . Write  $t := w_j$  and, for each input  $x$ , set  $z(t, x) := \mathbf{w}^\top x = t x_j + c_x$  with  $c_x := w_{-j}^\top x_{-j}$ . Throughout,  $\phi(z) = \max(0, z)$  and  $\|x\|_\infty \leq 1$ .

### 1. Derivative of $\Sigma$ and $J_A$

For  $\Sigma(\mathbf{w}) = \mathbb{E}_x[\phi(z)^2] = \mathbb{E}_x[z(t, x)^2 \mathbf{1}\{z(t, x) > 0\}]$  we have, for a.e.  $t$ ,

$$\partial_t \Sigma = 2 \mathbb{E}_x[z \mathbf{1}\{z > 0\} x_j], \quad \partial_t^2 \Sigma = 2 \mathbb{E}_x[x_j^2 \mathbf{1}\{z > 0\}] \leq 2, \quad (98)$$

because  $|x_j| \leq 1$ . By Cauchy–Schwarz,

$$|\partial_t \Sigma| \leq 2(\mathbb{E}_x[z^2 \mathbf{1}\{z > 0\}])^{1/2} (\mathbb{E}_x[x_j^2 \mathbf{1}\{z > 0\}])^{1/2} \leq 2\sqrt{\Sigma}. \quad (99)$$

For  $J_A(\mathbf{w}) = \mathbb{E}_x[\phi(z) \chi_A(x)]$  (with  $|\chi_A| \leq 1$ ), we get for a.e.  $t$ :

$$\partial_t J_A = \mathbb{E}_x[\mathbf{1}\{z > 0\} x_j \chi_A(x)], \quad \partial_t^2 J_A = 0 \text{ (a.e.)}, \quad |\partial_t J_A| \leq \mathbb{E}|x_j| \leq 1, \quad (100)$$

and by Cauchy–Schwarz again

$$|J_A| \leq (\mathbb{E}_x[\phi(z)^2])^{1/2} (\mathbb{E}_x[\chi_A(x)^2])^{1/2} = \sqrt{\Sigma}. \quad (101)$$

### 2. Decompose $g$ and bound each second derivative

Let

$$a := \sigma_a^{-2}, \quad c := \frac{1}{\kappa^2 N^{2\gamma}}, \quad J_\Delta := J_y - m_S J_S, \quad q := J_\Delta^2, \quad (102)$$

so

$$g(\mathbf{w}) = \underbrace{\frac{1}{2} \log(a + c\Sigma)}_{=: g_1(\Sigma)} - \underbrace{\frac{q}{2\kappa^4 N^{2\gamma} (a + c\Sigma)}}_{=: g_2(\Sigma, J_\Delta)}. \quad (103)$$

*Term 1:* Set  $h_1(s) := \frac{1}{2} \log(a + cs)$ , so  $h_1'(s) = \frac{c}{2(a+cs)}$  and  $h_1''(s) = -\frac{c^2}{2(a+cs)^2}$ . By the chain rule,

$$\partial_t^2 g_1 = h_1''(\Sigma) (\partial_t \Sigma)^2 + h_1'(\Sigma) \partial_t^2 \Sigma. \quad (104)$$

Using (99) and  $\partial_t^2 \Sigma \leq 2$ ,

$$\left| h_1''(\Sigma) (\partial_t \Sigma)^2 \right| \leq \frac{c^2}{2(a + c\Sigma)^2} \cdot 4\Sigma = 2c^2 \frac{\Sigma}{(a + c\Sigma)^2} \leq \frac{c\sigma_a^2}{2}, \quad (105)$$

where the last inequality follows by maximizing  $u \mapsto \frac{u}{(a+cu)^2}$  at  $u = a/c$ . Also  $|h_1'(\Sigma) \partial_t^2 \Sigma| \leq \frac{c}{2(a+c\Sigma)} \cdot 2 \leq c\sigma_a^2$ . Hence

$$|\partial_t^2 g_1| \leq \frac{3}{2} c\sigma_a^2. \quad (106)$$

*Term 2:* Write  $g_2 = -\alpha \frac{q}{a+c\Sigma}$  with  $\alpha := \frac{1}{2\kappa^4 N^{2\gamma}}$ . Differentiating twice and grouping terms gives (for a.e.  $t$ ):

$$\partial_t^2 g_2 = -\alpha \left[ \frac{q''}{a+c\Sigma} - \frac{2q'c\Sigma'}{(a+c\Sigma)^2} - \frac{qc\Sigma''}{(a+c\Sigma)^2} + \frac{2q(c\Sigma')^2}{(a+c\Sigma)^3} \right],$$

where primes denote  $\partial_t$ . We now bound  $q, q', q''$  with step 1 and using (101) and  $|\partial_t J_A| \leq 1$ .

- $q$ : Using  $|J_\Delta| \leq |J_Y| + |m_S| |J_S| \leq (1 + |m_S|)\sqrt{\Sigma} =: C_\Delta \sqrt{\Sigma}$  we get  $q = J_\Delta^2 \leq C_\Delta^2 \Sigma$ .
- $|q'|$ : We get

$$\begin{aligned} |q'| &= 2|J_\Delta| |\partial_t J_\Delta| \leq 2C_\Delta \sqrt{\Sigma} \cdot (|\partial_t J_Y| + |m_S| |\partial_t J_S|) \leq 2C_\Delta \sqrt{\Sigma} (1 + |m_S|) \\ &\leq 2C_\Delta^2 \sqrt{\Sigma} \end{aligned}$$

- $|q''|$ : We get  $q'' = 2(\partial_t J_\Delta)^2 \leq 2C_\Delta^2$

Together with  $\Sigma' \leq 2\sqrt{\Sigma}$  and  $\Sigma'' \leq 2$ , we get

$$\left| \frac{q''}{a+c\Sigma} \right| \leq \frac{2C_\Delta^2}{a}, \quad (107)$$

$$\left| \frac{2q'c\Sigma'}{(a+c\Sigma)^2} \right| \leq \frac{2 \cdot 2C_\Delta^2 \sqrt{\Sigma} \cdot c \cdot 2\sqrt{\Sigma}}{(a+c\Sigma)^2} = 8C_\Delta^2 c \frac{\Sigma}{(a+c\Sigma)^2} \leq \frac{2C_\Delta^2}{a}, \quad (108)$$

$$\left| \frac{qc\Sigma''}{(a+c\Sigma)^2} \right| \leq \frac{C_\Delta^2 \Sigma \cdot c \cdot 2}{(a+c\Sigma)^2} \leq \frac{C_\Delta^2}{a}, \quad (109)$$

$$\left| \frac{2q(c\Sigma')^2}{(a+c\Sigma)^3} \right| \leq \frac{2C_\Delta^2 \Sigma \cdot c^2 \cdot 4\Sigma}{(a+c\Sigma)^3} = 8C_\Delta^2 c^2 \frac{\Sigma^2}{(a+c\Sigma)^3} \leq \frac{32C_\Delta^2}{27a}, \quad (110)$$

where in the last two lines we used that  $u \mapsto \frac{u}{(a+cu)^2}$  and  $u \mapsto \frac{u^2}{(a+cu)^3}$  are maximized at  $u = \frac{a}{c}$  and  $u = \frac{2a}{c}$ , respectively, with finite maxima depending only on  $a, c$ . Therefore  $|\partial_t^2 g_2| \leq \alpha K_2$  for a constant  $K_2$  depending only on  $(a, c, m_S)$ , and independent of  $d, w_{-j}$  and  $t$ .

### 3: Uniform bound on $\partial_t^2 g$

Combining (106) and the bound for  $\partial_t^2 g_2$ , there exists

$$L_\star := \frac{3}{2} c \sigma_a^2 + \alpha K_2 \quad (111)$$

such that  $|\partial_t^2 g(w_{-j}, t)| \leq L_\star$  for a.e.  $t \in \mathbb{R}$ . This proves the first claim.

### 4: Standard Taylor expansion

For any twice-differentiable  $h$  with  $|h''| \leq L_\star$  a.e., the 1D Taylor inequality gives

$$h(t) \geq h(0) + h'(0)t - \frac{L_\star}{2} t^2 \quad (\forall t \in \mathbb{R}).$$

Applying this to  $h(t) = g(w_{-j}, t)$  yields

$$g(w_{-j}, t) \geq g(w_{-j}, 0) + \partial_t g(w_{-j}, 0)t - \frac{L_\star}{2} t^2.$$

In the parity setting and for  $j \notin S$  (off-support), symmetry implies  $\partial_t g(w_{-j}, 0) = 0$ , giving the stated global quadratic lower bound  $g(w_{-j}, t) \geq g(w_{-j}, 0) - \frac{L_\star}{2} t^2$ .  $\blacksquare$

### E.6. Proof of Theorem 2

**Proof** In plain MF, the single-neuron posterior has the form

$$p_{\text{MF}}(\mathbf{w} \mid m_S) \propto \exp\left(-\frac{\rho}{2}\|\mathbf{w}\|_2^2 - g(\mathbf{w})\right), \quad \rho := \frac{d}{\sigma_w^2},$$

with  $g$  as in Lemma 12. For the parity target  $y(x) = \chi_S(x)$ , the law is invariant under permutations within  $S$  and within  $S^c$ , hence the second moments are constant on each group:

$$q_{\text{on}} := \mathbb{E}[w_j^2] \ (j \in S), \quad q_{\text{off}} := \mathbb{E}[w_j^2] \ (j \notin S).$$

**Off-support.** Fix  $j \notin S$  and condition on  $w_{-j}$ . In plain MF we have  $\rho_j = \rho = \Omega(d)$ , and by off-support symmetry  $\partial_j g(w_{-j}, 0) = 0$ . Lemma ?? therefore yields

$$\mathbb{E}[w_j^2 \mid w_{-j}] \leq \frac{1 + o_d(1)}{\rho} \quad \text{uniformly in } w_{-j}.$$

Averaging over  $w_{-j}$  gives  $q_{\text{off}} \leq \frac{1+o_d(1)}{\rho}$ .

**On-support.** For  $j \in S$ , the same curvature bound  $|\partial_t^2 g(w_{-j}, t)| \leq L_\star$  implies the 1D conditional potential  $U_j(t) = \frac{\rho}{2}t^2 + g(w_{-j}, t)$  is  $(\rho - L_\star)$ -strongly convex, hence  $(w_j \mid w_{-j}) \leq (\rho - L_\star)^{-1}$ . Moreover  $|\partial_t g(w_{-j}, 0)|$  is  $O(1)$  uniformly in  $w_{-j}$  (by the same bounds used in Lemma 12), so the conditional mean satisfies  $|\mathbb{E}[w_j \mid w_{-j}]| = O(\rho^{-1})$ , and therefore

$$\mathbb{E}[w_j^2 \mid w_{-j}] = (w_j \mid w_{-j}) + \mathbb{E}[w_j \mid w_{-j}]^2 \leq \frac{1}{\rho - L_\star} + O(\rho^{-2}) = \frac{1 + O(1/\rho)}{\rho}.$$

Averaging gives  $q_{\text{on}} \leq \frac{1+O(1/\rho)}{\rho}$ .

Finally, the same argument applied to  $-g$  (or equivalently using the opposite quadratic envelope from  $|\partial_t^2 g| \leq L_\star$ ) gives matching lower bounds  $q_{\text{on}}, q_{\text{off}} \geq \frac{1-O(1/\rho)}{\rho}$ , so

$$\frac{q_{\text{on}}}{q_{\text{off}}} = 1 + O\left(\frac{1}{\rho}\right) = 1 + O\left(\frac{1}{d}\right).$$

Since  $\|\mathbf{w}_S\|_2^2 = \sum_{j \in S} w_j^2$  and  $\|\mathbf{w}_{S^c}\|_2^2 = \sum_{j \notin S} w_j^2$ , we get

$$R_{\mathbf{w}}^{\text{MF}} = \frac{\sqrt{d-k}\|\mathbf{w}_S\|_2}{\sqrt{k}\|\mathbf{w}_{S^c}\|_2} = \sqrt{\frac{q_{\text{on}}}{q_{\text{off}}}} = 1 + O\left(\frac{1}{d}\right).$$

■

### E.7. Proof of the moment closure

**Proof** At stationarity, the exact row-energy balance reads

$$\gamma \langle q_j \rangle_\infty + \langle \hat{\rho}_j q_j \rangle_\infty = TN.$$

Applying the linear response assumption

$$(\gamma + \rho_j^r) \langle q_j \rangle_\infty \approx TN.$$

Dividing by  $T$  and using the definition of  $\rho_{\text{eff},j}$  yields

$$\langle \mathbf{w}_j^2 \rangle_\infty \approx \frac{N}{\rho_{\text{eff},j}}.$$

Finally, by the MF assumption we have  $\langle \mathbf{w}_j^2 \rangle_\infty = Nw_j^2$ , hence  $w_j^2 \approx 1/\rho_{\text{eff},j}$ .

The MF-ARD fixed point equation for  $\rho_j$  is, by construction,

$$\rho_j = \frac{\alpha_0 + \frac{N}{2}}{\frac{\alpha_0}{d} + \frac{N}{2}Q_j},$$

For  $\alpha_0 \ll N$ , expanding the prefactor ratio gives  $\rho_j = (1/w_j^2)(1 + O(\alpha_0/N))$ , proving (??). This completes the proof upon identifying  $\rho_j \equiv \rho_{\text{eff},j}$ . ■

## Appendix G. Algorithms

### G.1. FP algorithm

**Model.** Given  $(X, y)$  with  $X \in \{\pm 1\}^{P \times d}$  and  $y \in \mathbb{R}^{P \times 1}$ , we approximate the predictor by a particle ensemble,

$$f(\mathbf{x}) = s_f \sum_{b=1}^B a_b \phi(\mathbf{w}_b^\top \mathbf{x}), \quad s_f = \frac{N^{1-\gamma}}{B}. \quad (112)$$

We draw a single dataset of size  $P$  once at initialization and keep it fixed for the entire run. On the training set let  $f_p = f(\mathbf{x}_p)$  and  $r_p = y_p - f_p$ . The Langevin temperature is fixed by the likelihood noise as  $T = 2\kappa^2$  (Section C.1). This choice makes SGLD asymptotically sample from the Bayesian posterior. Because all objectives and gradients are computed as empirical averages over this fixed sample (via  $1/P$  factors), the dynamics naturally exhibit finite- $P$  fluctuations.

**Sufficient statistics** We define the following low-rank statistics per particle  $b$ , with  $z_{pb} = \mathbf{x}_p^\top \mathbf{w}_b$ :

$$C_{1,b} = \sum_{p=1}^P \phi(z_{pb}) r_p, \quad C_{2,b} = \sum_{p=1}^P \phi(z_{pb})^2, \quad (113)$$

$$G_b^{\text{data}} = -\frac{2}{P} \sum_{p=1}^P \left( r_p - s_f a_b \phi(z_{pb}) \right) \phi'(z_{pb}) (s_f a_b) x_p \in \mathbb{R}^d. \quad (114)$$

These quantities are the only per-pass summaries we need to form gradients for  $a_b$  and  $w_b$ .

**Prior and SGLD potential** We impose a diagonal (ARD) Gaussian prior on the weights and an i.i.d. Gaussian prior on amplitudes:

$$E_{\text{prior}} = \frac{1}{2} \sum_{b=1}^B \rho^\top (\mathbf{w}_b \odot \mathbf{w}_b) + \frac{1}{2\sigma_a^2} \sum_{b=1}^B a_b^2, \quad \mathcal{L}_{\text{data}} = \frac{1}{P} \sum_{p=1}^P (y_p - f_p)^2, \quad (115)$$

and update parameters by Langevin dynamics on the potential

$$U(W, a) = T E_{\text{prior}} + \mathcal{L}_{\text{data}}. \quad (116)$$

**Gradients used by SGLD** From the streamed statistics we get closed-form gradients:

$$\nabla_{a_b} U = \frac{T}{\sigma_a^2} a_b - \frac{2s_f}{P} C_{1,b} + \frac{2s_f^2}{P} C_{2,b} a_b, \quad \nabla_{w_b} U = G_b^{\text{data}} + T(\rho \odot \mathbf{w}_b). \quad (117)$$

These are the only quantities used inside the inner SGLD loop.

**SGLD updates** We apply Euler–Maruyama steps with isotropic Gaussian noise:

$$w_b \leftarrow w_b - \eta \nabla_{w_b} U + \sqrt{2T\eta} \xi_{w_b}, \quad a_b \leftarrow a_b - \eta \nabla_{a_b} U + \sqrt{2T\eta} \xi_{a_b}, \quad (118)$$

where  $\xi_{w_b} \sim \mathcal{N}(0, I_d)$  and  $\xi_{a_b} \sim \mathcal{N}(0, 1)$ . We use polynomial decay on  $\eta$  always matching the SGLD-trained full NNs.

---

**Algorithm 1** Simple streaming SGLD for  $(a, w)$  with optional ARD
 

---

**Require:**  $(X, y)$ ;  $B, N, \gamma, \sigma_a, \sigma_w, \kappa$ ; activation  $\phi$ ; steps  $T_{\text{out}}$ , inner steps  $K$ , step size  $\eta$ ; optional ARD  $(\alpha_0, \beta_0, \lambda)$

**Ensure:** Final particles  $\{(w_b, a_b)\}_{b=1}^B$  and predictor  $f(x) = s_f \sum_b a_b \phi(w_b^\top x)$

1: **Init:**  $w_b \sim \mathcal{N}(0, \sigma_w^2 I_d/d)$ ,  $a_b \sim \mathcal{N}(0, \sigma_a^2)$ ; set  $\rho \leftarrow d/\sigma_w^2$ ; set  $T \leftarrow 2\kappa^2$  **for**  $t = 1..T_{\text{out}}$  **do**

2:

**end**

compute  $f_p = s_f \sum_b a_b \phi(x_p^\top w_b)$  and residuals  $r_p = y_p - f_p$  **for**  $k = 1..K$  **do**

3:

**end**

compute  $\{C_{1,b}, C_{2,b}, G_b^{\text{data}}\}$  via formulas above

4: Form  $\nabla_{a_b} U, \nabla_{w_b} U$ ; update  $w_b \leftarrow w_b - \eta \nabla_{w_b} U + \sqrt{2T\eta} \xi_{w_b}$ ,  $a_b \leftarrow a_b - \eta \nabla_{a_b} U + \sqrt{2T\eta} \xi_{a_b}$

5: refresh  $f_p, r_p$  after the last inner step

6:

7: ARD update  $\rho$ :  $\alpha_{\text{post}} = \alpha_0 + \frac{B}{2}$ ,  $\beta_{\text{post}} = \beta_0 + \frac{1}{2} \sum_b \|w_b\|^2$ ,  $\rho \leftarrow (1 - \lambda)\rho + \lambda \alpha_{\text{post}}/\beta_{\text{post}}$

8:

---

**ARD update** The ARD update is:

$$\alpha_{\text{post}} = \alpha_0 + \frac{B}{2}, \quad \beta_{\text{post}} = \beta_0 + \frac{1}{2} \sum_{b=1}^B \|w_b\|_2^2, \quad \rho \leftarrow (1 - \lambda)\rho + \lambda \frac{\alpha_{\text{post}}}{\beta_{\text{post}}}. \quad (119)$$

**Fixed-point view and the  $K$  inner steps** The outer loop implements a fixed-point iteration on the network  $f$ . Writing  $r = y - f$ , define the map  $\mathcal{G}_K$  as: (i) run  $K$  inner SGLD steps on  $(w_b, a)$  using the current residual  $r$ ; (ii) recompute  $f^{\text{new}}(x_p) = s_f \sum_b a_b \phi(w_b^\top x_p)$ . As  $K \rightarrow \infty$  and  $\eta \rightarrow 0$  the inner Markov chain approaches its stationary law, and the iteration solves the MF FP equations described in the theory sections.

**Empirical calculation of  $m_S$  and generalisation error** Let the teacher be a single Walsh mode  $\chi_S$ , so  $y(x) = \chi_S(x)$ . On a held-out set  $\{x_\mu\}_{\mu=1}^{P_{\text{eval}}}$  define the vector  $c \in \mathbb{R}^{P_{\text{eval}}}$  by  $c_\mu = \chi_S(x_\mu)$ , the empirical Gram scalar

$$g = \frac{1}{P_{\text{eval}}} c^\top c, \quad (120)$$

and the (scalar) empirical overlap

$$m_S = \frac{1}{P_{\text{eval}}} \sum_{\mu=1}^{P_{\text{eval}}} \chi_S(x_\mu) f(x_\mu) = \frac{1}{P_{\text{eval}}} c^\top f. \quad (121)$$

Let  $\bar{f}^2 = \frac{1}{P_{\text{eval}}} \sum_\mu f(x_\mu)^2$ . Then the empirical test MSE decomposes as

$$\hat{\mathcal{E}}_{\text{test}} = \frac{1}{2P_{\text{eval}}} \sum_{\mu=1}^{P_{\text{eval}}} (f(x_\mu) - \chi_S(x_\mu))^2 = \underbrace{\frac{1}{2}(1 - m_S)^2}_{\text{mode term}} + \underbrace{\frac{1}{2}(\bar{f}^2 - 2m_S^2 + g m_S^2)}_{\text{noise / orthogonal term}} \quad (122)$$

$$= \frac{1}{2} (\bar{f}^2 - 2m_S + g), \quad (123)$$

where the second line is the direct empirical expression. When the Walsh basis is orthonormal on the evaluation set ( $g = 1$ ), this reduces to  $\hat{\mathcal{E}}_{\text{test}} = \frac{1}{2}(\bar{f}^2 - 2m_S + 1)$ .

**Appendix H. Training details (hyperparameters)**

Here, we present the hyperparameters for SGLD and MF-ARD for the different figures. The hyperparameters for Figure 2 are specified below.

- Data in Figure 1 is a slice of Figure 2 for  $\kappa = 5 \cdot 10^{-3}$ .
- Data in Figure 4 is a slice of Figure 2 for  $\kappa = 5 \cdot 10^{-3}$ .

Hyperparameter	SGLD
$d$ (input dimension)	35
$P$ (train set sizes)	{10, 100, 500, 750, 1000, 2133, 3666, 5000, 7500, 10000}
$\kappa$ values	{ $5 \times 10^{-4}$ , $10^{-3}$ , $5 \times 10^{-3}$ , $7.5 \times 10^{-3}$ , $10^{-2}$ , $5 \times 10^{-2}$ , $10^{-1}$ }
$E$ (experiments / config)	3
teacher set $S$	{0, 1, 2, 3}
data distribution	$X \in \{-1, +1\}^d$ , $y = \prod_{j \in S} X_j$
activation	ReLU
$N$ (hidden units)	512
$\gamma$ (scaling exponent)	0.5
$g_w, g_a$ (prior variances)	1.0, 1.0
initialization	$w \sim \mathcal{N}(0, g_w/d)$ , $a \sim \mathcal{N}(0, g_a)$
temperature $T$	$2 \kappa^2$
epochs (max)	7,500,000
batch size	full-batch
loss	mean MSE
learning rate $\eta$ (final)	$5 \times 10^{-4}$
start LR $\eta_{\text{start}}$	$5 \times 10^{-3}$
LR scheduler	polynomial (power 2): $\eta_{\text{start}} \rightarrow \eta$ over $2 \times 10^6$ steps

Table 2: Algorithm outlined in Section C.1. Hyperparameters for Figure 2 a).

Hyperparameter	MF-ARD
$d$ (input dimension)	35
$P$ (train set sizes)	{10, 100, 500, 750, 1000, 2133, 3666, 5000, 7500, 10000}
$\kappa$ values	{ $5 \times 10^{-4}$ , $10^{-3}$ , $5 \times 10^{-3}$ , $7.5 \times 10^{-3}$ , $10^{-2}$ , $5 \times 10^{-2}$ , $10^{-1}$ }
$E$ (experiments / config)	3
teacher set $S$	{0, 1, 2, 3}
data distribution	$X \in \{-1, +1\}^d$ , $y = \prod_{j \in S} X_j$
activation	ReLU
$B$ (particles)	512
$N$	512
$\gamma$	0.5
$\sigma_a, \sigma_w$	1.0, 1.0
initialization	$w \sim \mathcal{N}(0, \sigma_w^2/d)$ , $a \sim \mathcal{N}(0, \sigma_a^2)$
outer steps	7,500,000
learning rate scheduler	poly-2 decay: $5 \times 10^{-3} \rightarrow 5 \times 10^{-4}$ over $2 \times 10^6$ steps
SGLD inner steps $K$	$K_0=12 \rightarrow K_{\min}=2$ (decay over $6 \times 10^5$ steps)
temperature $T$	$2\kappa^2$
ARD	on: $\alpha_0=4.0$ , $\beta_0=4/35$ , EMA 0.5, $\rho \in [0, 10^{18}]$

Table 3: Algorithm outlined in Section G.1. Hyperparameters for Figure 2 b) with ARD disabled and c) with ARD enabled.

For the single index model we introduced a bias vector in the architecture.

Hyperparameter	SGLD (Hermite single-index)
$d$ (input dimension)	18
$P$ (train set sizes)	{75,000, 50,000, 25,000, 10,000, 5,000, 1,000, 100, 50}
$\kappa$ values	$\{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$
$E$ (experiments / config)	4
teacher type	single-index Hermite
Hermite degree $p$	4
support size $k$	2 (random per experiment)
data distribution	$X \sim \mathcal{N}(0, I_d)$
teacher $w$	$w_i = \frac{1}{\sqrt{k}}$ on support, else 0
labels	$y = \text{He}_p(Xw)$
activation	ReLU
$N$	1024
$\gamma$	0.5
$\sigma_a, \sigma_w, \sigma_b$	1.0, 0.5, 1.0
initialization	$w \sim \mathcal{N}(0, \sigma_w^2/d), a \sim \mathcal{N}(0, \sigma_a^2), b \sim \mathcal{N}(0, \sigma_b^2)$
temperature $T$	$2\kappa^2$
epochs (max)	4,000,000
batch size	full-batch
loss	mean MSE
learning rate $\eta$ (final)	$5 \times 10^{-4}$
start LR $\eta_{\text{start}}$	$2 \times 10^{-3}$
LR scheduler	polynomial (power 2): $\eta_{\text{start}} \rightarrow \eta$ over $2 \times 10^6$ steps

Table 4: Algorithm outlined in Section C.1. Hyperparameters for Figure 8 a).

Hyperparameter	MF-ARD (Hermite single-index)
$d$ (input dimension)	18
$P$ (train set sizes)	{75,000, 50,000, 25,000, 10,000, 5,000, 1,000, 100, 50}
$\kappa$ values	$\{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$
$E$ (experiments / config)	4
teacher type	single-index Hermite
Hermite degree $p$	4
support size $k$	2 (random per experiment)
data distribution	$X \sim \mathcal{N}(0, I_d)$
teacher $w$	$w_i = \frac{1}{\sqrt{k}}$ on support, else 0
labels	$y = \text{He}_p(Xw)$
activation	ReLU
$B$ (particles)	1024
$N$	1024
$\gamma$	0.5
$\sigma_a, \sigma_w, \sigma_b$	1.0, 0.5, 1.0
initialization	$w \sim \mathcal{N}(0, \sigma_w^2/d), a \sim \mathcal{N}(0, \sigma_a^2), b \sim \mathcal{N}(0, \sigma_b^2)$
outer steps	4,000,000
learning rate scheduler	poly-2 decay: $2 \times 10^{-3} \rightarrow 5 \times 10^{-4}$ over $2.5 \times 10^6$ steps
SGLD inner steps $K$	$K_0=12 \rightarrow K_{\min}=2$ (decay over $6 \times 10^5$ steps)
temperature $T$	$2\kappa^2$
ARD	on: $\alpha_0=0.1, \beta_0=\alpha_0/d=0.1/18, \text{EMA } 0.5, \rho \in [0, 10^{18}]$

Table 5: Algorithm outlined in Section G.1. Hyperparameters for Figure 8 b) with ARD enabled.