
Dynamics of Finite Width Kernel and Prediction Fluctuations in Mean Field Neural Networks

Blake Bordelon & Cengiz Pehlevan

John Paulson School of Engineering and Applied Sciences,
Center for Brain Science,
Kempner Institute for the Study of Natural & Artificial Intelligence,
Harvard University
Cambridge MA, 02138
blake_bordelon@g.harvard.edu, cpehlevan@g.harvard.edu

Abstract

We analyze the dynamics of finite width effects in wide but finite feature learning neural networks. Starting from a dynamical mean field theory description of infinite width deep neural network kernel and prediction dynamics, we provide a characterization of the $\mathcal{O}(1/\sqrt{\text{width}})$ fluctuations of the DMFT order parameters over random initializations of the network weights. Our results, while perturbative in width, unlike prior analyses, are non-perturbative in the strength of feature learning. In the lazy limit of network training, all kernels are random but static in time and the prediction variance has a universal form. However, in the rich, feature learning regime, the fluctuations of the kernels and predictions are dynamically coupled with a variance that can be computed self-consistently. In two layer networks, we show how feature learning can dynamically reduce the variance of the final tangent kernel and final network predictions. We also show how initialization variance can slow down online learning in wide but finite networks. In deeper networks, kernel variance can dramatically accumulate through subsequent layers at large feature learning strengths, but feature learning continues to improve the signal-to-noise ratio of the feature kernels. In discrete time, we demonstrate that large learning rate phenomena such as edge of stability effects can be well captured by infinite width dynamics and that initialization variance can decrease dynamically. For CNNs trained on CIFAR-10, we empirically find significant corrections to both the bias and variance of network dynamics due to finite width.

1 Introduction

Learning dynamics of deep neural networks are challenging to analyze and understand theoretically, but recent progress has been made by studying the idealization of infinite-width networks. Two types of infinite-width limits have been especially fruitful. First, the kernel or lazy infinite-width limit, which arises in the standard or neural tangent kernel (NTK) parameterization, gives prediction dynamics which correspond to a linear model [1–5]. This limit is theoretically tractable but fails to capture adaptation of internal features in the neural network, which are thought to be crucial to the success of deep learning in practice. Alternatively, the mean field or μ -parameterization allows feature learning at infinite width [6–9].

With a set of well-defined infinite-width limits, prior theoretical works have analyzed finite networks in the NTK parameterization perturbatively, revealing that finite width both enhances the amount of feature evolution (which is still small in this limit) but also introduces variance in the kernels and the predictions over random initializations [10–15]. Because of these competing effects, in some situations wider networks are better, and in others wider networks perform worse [16].

In this paper, we analyze finite-width network learning dynamics in the mean field parameterization. In this parameterization, wide networks are empirically observed to outperform narrow networks [7, 17, 18]. Our results and framework provide a methodology for reasoning about detrimental finite-size effects in such feature-learning neural networks. We show that observable averages involving kernels and predictions obey a well-defined power series in inverse width even in rich training regimes. We generally observe that the leading finite-size corrections to both the bias and variance components of the square loss are increased for narrower networks, and diminish performance. Further, we show that richer networks are closer to their corresponding infinite-width mean field limit. For simple tasks and architectures the leading $\mathcal{O}(1/\text{width})$ corrections to the error can be descriptive, while for large sample size or more realistic tasks, higher order corrections appear to become relevant. Concretely, our contributions are listed below:

1. Starting from a dynamical mean field theory (DMFT) description of infinite-width nonlinear deep neural network training dynamics, we provide a complete recipe for computing fluctuation dynamics of DMFT order parameters over random network initializations during training. These include the variance of the training and test predictions and the $\mathcal{O}(1/\text{width})$ variance of feature and gradient kernels throughout training.
2. We first solve these equations for the lazy limit, where no feature learning occurs, recovering a simple differential equation which describes how prediction variance evolves during learning.
3. We solve for variance in the rich feature learning regime in two-layer networks and deep linear networks. We show richer nonlinear dynamics improve the signal-to-noise ratio (SNR) of kernels and predictions, leading to closer agreement with infinite-width mean field behavior.
4. We analyze in a two-layer model why larger training set sizes in the overparameterized regime enhance finite-width effects and how richer training can reduce this effect.
5. We show that large learning rate effects such as edge-of-stability [19–21] dynamics can be well captured by infinite width theory, with finite size variance accurately predicted by our theory.
6. We test our predictions in Convolutional Neural Networks (CNNs) trained on CIFAR-10 [22]. We observe that wider networks and richly trained networks have lower logit variance as predicted. However, the timescale of training dynamics is significantly altered by finite width even after ensembling. We argue that this is due to a detrimental correction to the mean dynamical NTK.

1.1 Related Works

Infinite-width networks at initialization converge to a Gaussian process with a covariance kernel that is computed with a layerwise recursion [23–26, 13]. In the large but finite width limit, these kernels do not concentrate at each layer, but rather propagate finite-size corrections forward through the network [27–30, 14]. During gradient-based training with the NTK parameterization, a hierarchy of differential equations have been utilized to compute small feature learning corrections to the kernel through training [10–13]. However the higher order tensors required to compute the theory are initialization dependent, and the theory breaks down for sufficiently rich feature learning dynamics. Various works on Bayesian deep networks have also considered fluctuations and perturbations in the kernels at finite width during inference [31, 32]. Other relevant work in this domain are [33–39].

An alternative to standard/NTK parameterization is the mean field or μP -limit where features evolve even at infinite width [6–9, 40–42]. Recent studies on two-layer mean field networks trained online with Gaussian data have revealed that finite networks have larger sensitivity to SGD noise [43, 44]. Here, we examine how finite-width neural networks are sensitive to initialization noise. Prior work has studied how the weight space distribution and predictions converge to mean field dynamics with a dynamical error $\mathcal{O}(1/\sqrt{\text{width}})$ [40, 45], however in the deep case this requires a probability distribution over couplings between adjacent layers. Our analysis, by contrast, focuses on a function and kernel space picture which decouples interactions between layers at infinite width. A starting point for our present analysis of finite-width effects was a previous set of studies [9, 46] which identified the DMFT action corresponding to randomly initialized deep NNs which generates the distribution over kernel and network prediction dynamics. These prior works discuss the possibility of using a finite-size perturbation series but crucially failed to recognize the role of the network prediction fluctuations on the kernel fluctuations which are necessary to close the self-consistent equations in the rich regime. Using the mean field action to calculate a perturbation expansion around DMFT is a long celebrated technique to obtain finite size corrections in physics [47–50] and has been utilized for random untrained recurrent networks [51, 52], and more recently to calculate variance of

feature kernels Φ^\cdot at initialization $t = 0$ in deep MLPs or RNNs [53]. We extend these prior studies to the dynamics of training and to probe how feature learning alters finite size corrections.

2 Problem Setup

We consider wide neural networks where the number of neurons (or channels for a CNN) N in each layer is large. For a multi-layer perceptron (MLP), the network is defined as a map from input $\mathbf{x} \in \mathbb{R}^D$ to hidden preactivations $\mathbf{h} \in \mathbb{R}^N$ in layers $\ell \in \{1, \dots, L\}$ and finally output f

$$f = \frac{1}{\gamma N} \mathbf{w}^L \cdot \phi(\mathbf{h}^L), \quad \mathbf{h}^{\ell+1} = \frac{1}{\sqrt{N}} \mathbf{W}^\ell \phi(\mathbf{h}^\ell), \quad \mathbf{h}^1 = \frac{1}{\sqrt{D}} \mathbf{W}^0 \mathbf{x}, \quad (1)$$

where γ is a scale factor that controls feature learning strength, with large γ leading to rich feature learning dynamics and the limit of small $\gamma \rightarrow 0$ (or generally if γ scales as $N^{-\alpha}$ for $\alpha > 0$ as $N \rightarrow \infty$, NTK parameterization corresponds to $\alpha = \frac{1}{2}$) gives lazy learning where no features are learned [4, 7, 9]. The parameters $\mathbf{W} = \{\mathbf{W}^0, \mathbf{W}^1, \dots, \mathbf{W}^L\}$ are optimized with gradient descent or gradient flow $\frac{d}{dt} = -N\gamma^2 \nabla \mathcal{L}$ where $\mathcal{L} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \ell(f(\mathbf{x}), y)$ is a loss computed over dataset $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_P, y_P)\}$. This parameterization and learning rate scaling ensures that $\frac{d}{dt} f \sim \mathcal{O}_N$ (1) and $\frac{d}{dt} \mathbf{h}^\ell \sim \mathcal{O}_N$ (γ) at initialization. This is equivalent to maximal update parameterization (μ P)[8], which can be easily extended to other architectures including neural networks with trainable bias parameters as well as convolutional, recurrent, and attention layers [8, 9].

3 Review of Dynamical Mean Field Theory

The infinite-width training dynamics of feature learning neural networks was described by a DMFT in [9, 46]. We first review the DMFT's key concepts, before extending it to get insight into finite-widths. To arrive at the DMFT, one first notes that the training dynamics of such networks can be rewritten in terms of a collection of dynamical variables (or *order parameters*) $\mathbf{q} = \text{Vec}\{f^\cdot(t), \Phi^\cdot(t, s), G^\cdot(t, s), \dots\}$ [9], which include feature and gradient kernels [9, 54]

$$\Phi^\cdot(t, s) \equiv \frac{1}{N} \phi(\mathbf{h}^\cdot(t)) \cdot \phi(\mathbf{h}^\cdot(s)), \quad G^\cdot(t, s) \equiv \frac{1}{N} \mathbf{g}^\cdot(t) \cdot \mathbf{g}^\cdot(s), \quad (2)$$

where $\mathbf{g}^\cdot(t) = \gamma N \frac{\partial f^\cdot(t)}{\partial \mathbf{h}^\cdot(t)}$ are the back-propagated gradient signals. Further, for width- N networks the distribution of these dynamical variables across weight initializations (from a Gaussian distribution $\sim \mathcal{N}(0, \mathbf{I})$) is given by $p(\mathbf{q}) \propto \exp(N S(\mathbf{q}))$, where the action $S(\mathbf{q})$ contains interactions between neuron activations and the kernels at each layer [9].

The DMFT introduced in [9] arises in the $N \rightarrow \infty$ limit when $p(\mathbf{q})$ is strongly peaked around the saddle point \mathbf{q}_1 where $\frac{\partial S}{\partial \mathbf{q}}|_{\mathbf{q}_1} = 0$. Analysis of the saddle point equations reveal that the training dynamics of the neural network can be alternatively described by a stochastic process. A key feature of this process is that it describes the training time evolution of the distribution of neuron pre-activations in each layer (informally the histogram of the elements of $\mathbf{h}^\cdot(t)$) where each neuron's pre-activation behaves as an i.i.d. draw from this *single-site* stochastic process. We denote these random processes by $h^\cdot(t)$. Kernels in (2) are now computed as *averages* over these infinite-width single site processes $\Phi^\cdot(t, s) = \overline{\phi(h^\cdot(t)) \phi(h^\cdot(s))}$, $G^\cdot(t, s) = \overline{\mathbf{g}^\cdot(t) \mathbf{g}^\cdot(s)}$, where averages arise from the $N \rightarrow \infty$ limit of the dot products in (2). DMFT also provides a set of self-consistent equations that describe the complete statistics of these random processes, which depend on the kernels, as well as other quantities. To make our notation and terminology clearer for a machine learning audience, we provide Table 1 for a definition of the physics terminology in machine learning language.

Order params. \mathbf{q}	Action $S(\mathbf{q})$	Propagator	Single Site Density
Concentrating variables	\mathbf{q} 's log-density	Asymptotic Covariance	Neuron Marginals

Table 1: Relationship between the physics and ML terminology for the central objects in this paper. The \mathbf{q} which concentrate at infinite width, but fluctuate at finite width N . This paper is primarily interested in \mathbf{q} , the asymptotic covariance of the order parameters.

4 Dynamical Fluctuations Around Mean Field Theory

We are interested in going beyond the infinite-width limit to study more realistic finite-width networks. In this regime, the order parameters \mathbf{q} fluctuate in a $\mathcal{O}(1/\sqrt{N})$ neighborhood of \mathbf{q}_1 [55, 51, 53, 46].

Statistics of these fluctuations can be calculated from a general cumulant expansion (see App. D) [55, 56, 51]. We will focus on the leading-order corrections to the infinite-width limit in this expansion.

Proposition 1 *The finite-width N average of observable $O(\mathbf{q})$ across initializations, which we denote by $\langle O(\mathbf{q}) \rangle_N$, admits an expansion of the form whose leading terms are*

$$\langle O(\mathbf{q}) \rangle_N = \frac{\int d\mathbf{q} \exp(N S[\mathbf{q}]) O(\mathbf{q})}{\int d\mathbf{q} \exp(N S[\mathbf{q}])} = \langle O(\mathbf{q}) \rangle_1 + N [\langle V(\mathbf{q}) O(\mathbf{q}) \rangle_1 - \langle V(\mathbf{q}) \rangle_1 \langle O(\mathbf{q}) \rangle_1] + \dots, \quad (3)$$

where $\langle \cdot \rangle_1$ denotes an average over the Gaussian distribution $\mathbf{q} \sim \mathcal{N}(\mathbf{q}_1, -\frac{1}{N} \nabla_{\mathbf{q}}^2 S[\mathbf{q}_1])^{-1}$ and the function $V(\mathbf{q}) \equiv S(\mathbf{q}) - S(\mathbf{q}_1) - \frac{1}{2}(\mathbf{q} - \mathbf{q}_1)^T \nabla_{\mathbf{q}}^2 S(\mathbf{q}_1)(\mathbf{q} - \mathbf{q}_1)$ contains cubic and higher terms in the Taylor expansion of S around \mathbf{q}_1 . The terms shown include all the leading and sub-leading terms in the series in powers of $1/N$. The terms in ellipses are at least $\mathcal{O}(N^{-1})$ suppressed compared to the terms provided.

The proof of this statement is given in App. D. The central object to characterize finite size effects is the unperturbed covariance (the *propagator*): $\mathbb{D} = -\nabla_{\mathbf{q}}^2 S(\mathbf{q}_1) \mathbb{D}^{-1}$. This object can be shown to capture leading order fluctuation statistics $\langle (\mathbf{q} - \mathbf{q}_1)(\mathbf{q} - \mathbf{q}_1)^T \rangle_N = \frac{1}{N} \mathbb{D} + \mathcal{O}(N^{-2})$ (App. D.1), which can be used to reason about, for example, expected square error over random initializations. Correction terms at finite width may give a possible explanation of the superior performance of wide networks at fixed γ [7, 17, 18]. To calculate such corrections, in App. E, we provide a complete description of Hessian $\nabla_{\mathbf{q}}^2 S(\mathbf{q})$ and its inverse (the propagator) for a depth- L network. This description constitutes one of our main results. The resulting expressions are lengthy and are left to App. E. Here, we discuss them at a high level. Conceptually there are two primary ingredients for obtaining the full propagator:

- Hessian sub-blocks κ which describe the *uncoupled variances* of the kernels, such as

$$\kappa(t, s, t^\ell, s^\ell) \equiv \phi(h(t))\phi(h(s))\phi(h(t^\ell))\phi(h(s^\ell)) - \Phi(t, s)\Phi(t^\ell, s^\ell) \quad (4)$$

Similar terms also appear in other studies on finite width Bayesian inference [13, 31, 32] and in studies on kernel variance at initialization [27, 14, 29, 53].

- Blocks which capture the *sensitivity* of field averages to perturbations of order parameters, such as

$$D^{-1}(t, s, t^\ell, s^\ell) \equiv \frac{\partial \phi(h(t))\phi(h(s))}{\partial \Phi^{-1}(t^\ell, s^\ell)}, \quad D^G(t, s, t^\ell) \equiv \frac{\partial g(t)g(s)}{\partial \Delta(t^\ell)}, \quad (5)$$

where $\Delta(t) = -\frac{\partial \langle f \cdot y \rangle}{\partial f} \Big|_{f(t)}$ are error signal for each data point.

Abstractly, we can consider the uncoupled variances as “sources” of finite-width noise for each order parameter and the D blocks as summarizing a directed causal graph which captures how this noise propagates in the network (through layers and network predictions). In Figure 1, we illustrate this graph showing directed lines that represent causal influences of order parameters on fields and vice versa. For instance, if Φ were perturbed, D^{-1} would quantify the resulting perturbation to Φ^{-1} through the fields h^{-1} .

In App. E, we calculate \mathbb{D} and D tensors, and show how to use them to calculate the propagator. As an example of our results:

Proposition 2 *Partition \mathbf{q} into primal $\mathbf{q}_1 = \text{Vec}\{f(t), \Phi(t, s), \dots\}$ and conjugate variables $\mathbf{q}_2 = \text{Vec}\{\hat{f}(t), \hat{\Phi}(t, s), \dots\}$. Let $\mathbb{D} = \frac{\partial^2}{\partial \mathbf{q}_2 \partial \mathbf{q}_2} S[\mathbf{q}_1, \mathbf{q}_2]$ and $D = \frac{\partial^2}{\partial \mathbf{q}_2 \partial \mathbf{q}_1} S[\mathbf{q}_1, \mathbf{q}_2]$, then the propagator for \mathbf{q}_1 has the form $\mathbb{D}^{-1} = D^{-1} \mathbb{D}^{-1}$ (App E). The variables \mathbf{q}_1 are related to network observables, while conjugates \mathbf{q}_2 arise as Lagrange multipliers in the DMFT calculation. From the propagator \mathbb{D}^{-1} we can read off the variance of network observables such as $N \text{Var}(f) \sim \Sigma_f$.*

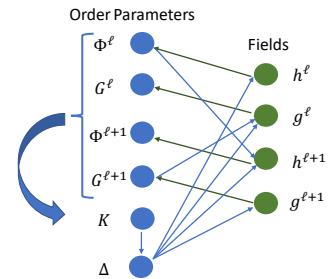


Figure 1: The directed causal graph between DMFT order parameters (blue) and fields (green) defines the D tensors of our theory. Each arrow represents a causal dependence. K denotes the NTK.

The necessary order parameters for calculating the fluctuations are obtained by solving the DMFT using numerical methods introduced in [9]. We provide a pseudocode for this procedure in App. F. We proceed to solve the equations defining in special cases which are illuminating and numerically feasible including lazy training, two layer networks and deep linear NNs.

5 Lazy Training Limit

To gain some initial intuition about why kernel fluctuations alter learning dynamics, we first analyze the static kernel limit $\gamma \rightarrow 0$ where features are frozen. To prevent divergence of the network in this limit, we use a background subtracted function $\tilde{f}(\mathbf{x}, \mathbf{x}^0) = f(\mathbf{x}, \mathbf{x}^0) - f(\mathbf{x}, \mathbf{x}^0)$ which is identically zero at initialization [4]. For mean square error, the $N \rightarrow \infty$ and $\gamma \rightarrow 0$ limit is governed by $\frac{\partial f(\mathbf{x})}{\partial t} = \mathbb{E}_{\mathbf{x}^0 \sim \mathcal{D}} \Delta(\mathbf{x}^0) K(\mathbf{x}, \mathbf{x}^0)$ with $\Delta(\mathbf{x}) = y(\mathbf{x}) - \tilde{f}(\mathbf{x})$ (for MSE) and K is the static (finite width and random) NTK. The finite- N initial covariance of the NTK has been analyzed in prior works [27, 13, 14], which reveal a dependence on depth and nonlinearity. Since the NTK is static in the $\gamma \rightarrow 0$ limit, it has constant initialization variance through training. Further, all sensitivity blocks of the Hessian involving the kernels and the prediction errors (such as the $D_{\mathbf{x}^0}$) vanish. We represent the covariance of the NTK as $\kappa(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4) = N \text{Cov}(K(\mathbf{x}_1, \mathbf{x}_2), K(\mathbf{x}_3, \mathbf{x}_4))$. To identify the dynamics of the error covariance, we relate K , the finite width NTK, to K_γ which is the deterministic infinite width NTK $K_\gamma(\mathbf{x}, \mathbf{x}^0) = \sum_k \lambda_k \psi_k(\mathbf{x}) \psi_k(\mathbf{x}^0)$ with respect to the training distribution \mathcal{D} , and decompose κ in this basis.

$$\kappa_{k \cdot mn} = \langle \kappa(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4) \psi_k(\mathbf{x}_1) \psi_n(\mathbf{x}_2) \psi_m(\mathbf{x}_3) \psi_m(\mathbf{x}_4) \rangle_{\mathbf{x}_1; \mathbf{x}_2; \mathbf{x}_3; \mathbf{x}_4 \sim \mathcal{D}}, \quad (6)$$

where averages are computed over the training distribution \mathcal{D} .

Proposition 3 For MSE loss, the prediction error covariance $\Sigma_k(t, s) = N \text{Cov}_0(\Delta_k(t), \Delta_k(s))$ satisfies a differential equation (App. H)

$$\frac{\partial}{\partial t} + \lambda_k \Sigma_k(t, s) = \frac{\partial}{\partial s} + \lambda \cdot \Sigma_k(t, s) = \sum_{nm} \kappa_{k \cdot mn} \Delta_m^\gamma(t) \Delta_n^\gamma(s), \quad (7)$$

where $\Delta_k^\gamma(t) \equiv \exp(-\lambda_k t) \langle \psi_k(\mathbf{x}) y(\mathbf{x}) \rangle_{\mathbf{x}}$ are the errors at infinite width for eigenmode k .

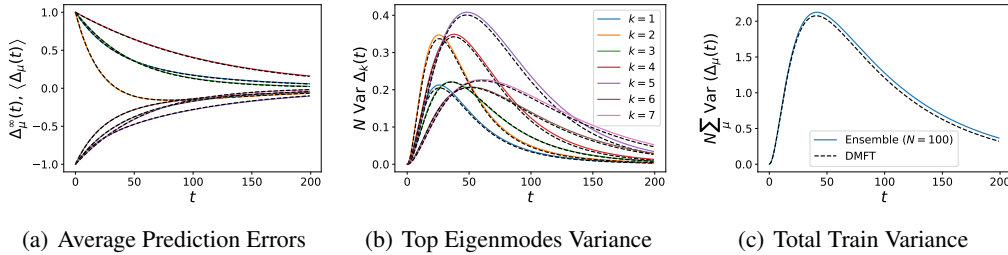


Figure 2: We show the accuracy of the lazy-limit ODE in equation (7) compared to a two-layer finite width $N = 100$ ReLU network trained with $\gamma = 0.05$ on $P = 10$ random training data points. (a) The average dynamics over an ensemble of $E = 500$ networks (solid colors) compared to the infinite width predictions (dashed black). (b) The predicted finite size variance for each eigenmode of the error $\Delta_k(t) = \sum_n \Sigma_k(t, s) \Delta_n^\gamma(s)$. These are not ordered simply by magnitude of eigenvalues or the target projections $y_k = \mathbf{y} \cdot \psi_k$, but rather depend on all eigenvalue gaps $\lambda_k - \lambda_\ell$ for $k \neq \ell$ and also the $\kappa_{k \cdot mn}$ tensor. (c) The total variance for all training points $N \sum_k \text{Var} \Delta_k(t) = N \sum_k \text{Var} \Delta_k(t)$ is also well predicted by the DMFT propagator equations.

An example verifying these dynamics is provided in Figure 2. In the case where the target is an eigenfunction $y = \psi_k$, the covariance has the form $\Sigma_k(t, s) = \kappa_{k \cdot k} \frac{\exp(-\lambda_k(t+s))}{(\lambda_k - \lambda_\ell)(\lambda_\ell - \lambda_k)}$. If the kernel is rank one with eigenvalue λ , then the dynamics have the simple form $\Sigma(t, s) = \kappa y^2 t s e^{-(t+s)}$. We note that similar terms appear in the prediction dynamics obtained by truncating the Neural Tangent Hierarchy [10, 11], however those dynamics concerned small feature learning

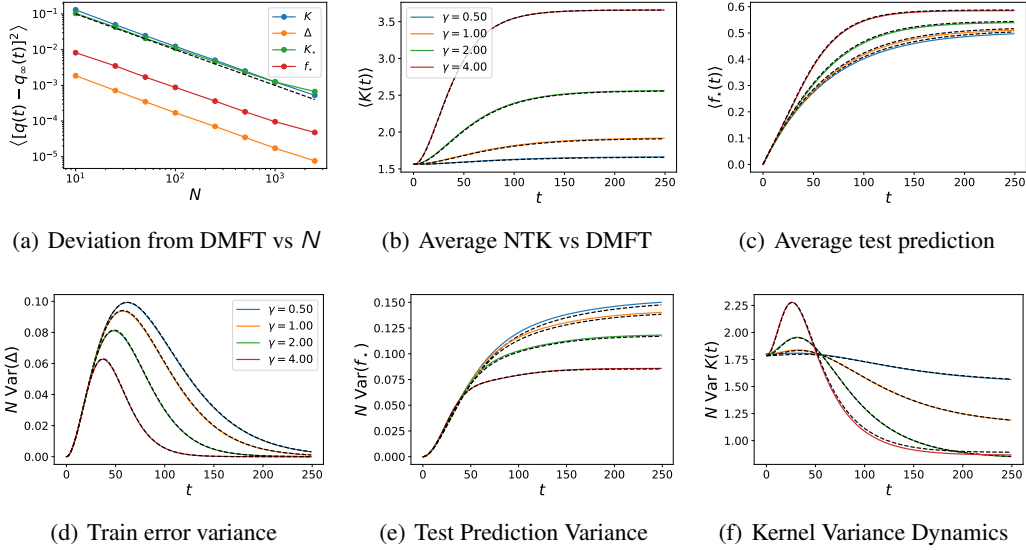


Figure 3: An ensemble of $E = 1000$ two layer $N = 256$ tanh networks trained on a single training point. Dashed black lines are DMFT predictions. (a) The square deviation from the infinite width DMFT scales as $\mathcal{O}(1/N)$ for all order parameters. (b) The ensemble average NTK $\langle K(t) \rangle$ (solid colors) and (c) ensemble average test point predictions $f_{\gamma}(t)$ for a point with $\frac{x \cdot x_{\gamma}}{D} = 0.5$ closely follow the infinite width predictions (dashed black). (d) The variance (estimated over the ensemble) of the train error $\Delta(t) = y - f(t)$ initially increases and then decreases as the training point is fit. (e) The variance of f_{γ} increases with time but decreases with γ . (f) The variance of the NTK during feature learning experiences a transient increase before decreasing to a lower value.

corrections rather than from initialization variance (App. H.1). Corrections to the mean $\langle \Delta \rangle$ are analyzed in App. H.2. We find that the variance and mean correction dynamics involves non-trivial coupling across eigendirections with a mixture of exponentials with timescales $\{\lambda_k^{-1}\}$.

6 Rich Regime in Two-Layer Networks

In this section, we analyze how feature learning alters the variance through training. We show a denoising effect where the signal to noise ratios of the order parameters improve with feature learning.

6.1 Kernel and Error Coupled Fluctuations on Single Training Example

In the rich regime, the kernel evolves over time but inherits fluctuations from the training errors. To gain insight, we first study a simplified setting where the data distribution is a single training example \mathbf{x} and single test point \mathbf{x}_{γ} in a two layer network. We will track $\Delta(t) = y - f(\mathbf{x}, t)$ and the test prediction $f_{\gamma}(t) = f(\mathbf{x}_{\gamma}, t)$. To identify the dynamics of these predictions we need the NTK $K(t)$ on the train point, as well as the train-test NTK $K_{\gamma}(t)$. In this case, all order parameters can be viewed as scalar functions of a single time index (unlike the deep network case, see App. E).

Proposition 4 *Computing the Hessian of the DMFT action and inverting (App. I), we obtain the following covariance for $\mathbf{q}_1 = \text{Vec}\{\Delta(t), f_{\gamma}(t), K(t), K_{\gamma}(t)\}_{t \in \mathbb{R}_+}$*

$$\mathbf{q}_1 = \begin{bmatrix} 2\mathbb{I} + \mathcal{K} & 0 & 0 & 0 \\ -\mathcal{K}_{\gamma} & \mathbb{I} & 0 & 0 \\ -\mathcal{D} & 0 & \mathbb{I} & 0 \\ -\mathcal{D}_{\gamma} & 0 & 0 & \mathbb{I} \end{bmatrix}^{-1} \quad (8)$$

where $[\mathcal{K}](t, s) = \Theta(t - s)K(s)$, $[\mathcal{K}_{\gamma}](t, s) = \Theta(t - s)K_{\gamma}(s)$ are Heaviside step functions and $D(t, s) = \frac{\partial^2}{\partial (s)}(\phi(h(t))^2 + g(t)^2)$ and $D_{\gamma}(t, s) = \frac{\partial^2}{\partial (s)}(\phi(h(t))\phi(h_{\gamma}(t)) + g(t)g_{\gamma}(t))$

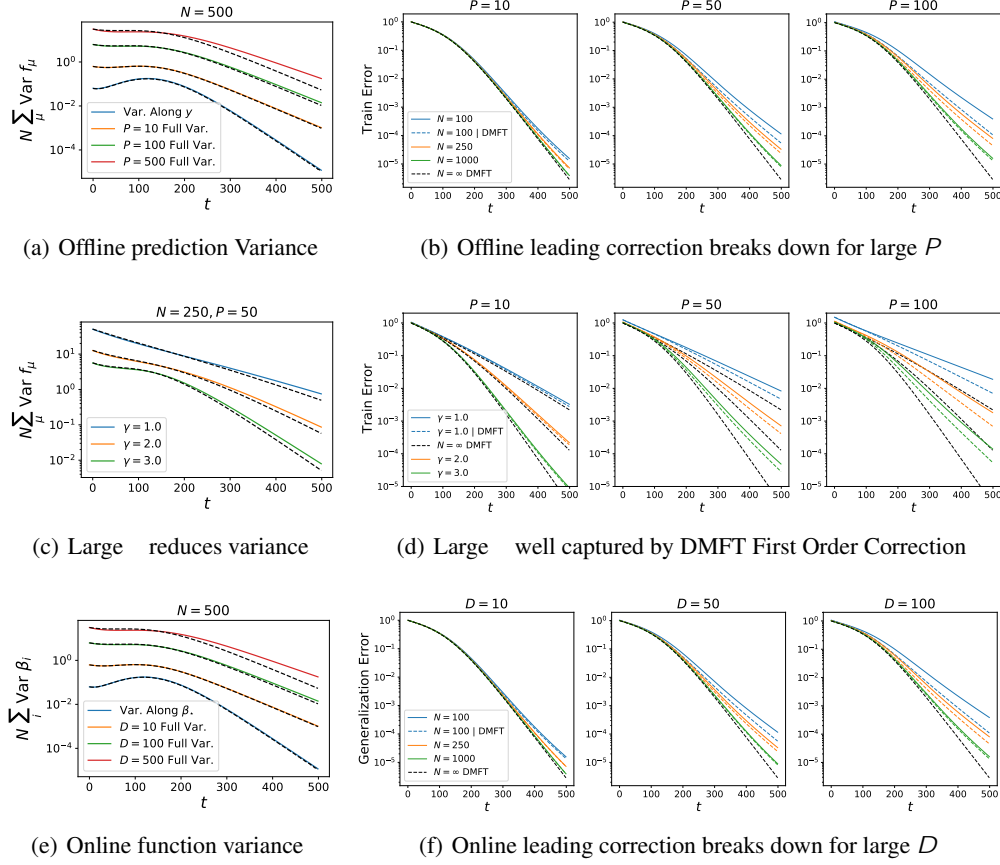


Figure 4: Large input dimension or multiple samples amplify finite size effects in a simple two layer model with unstructured data. Black dashed lines are theory. (a) The variance of offline learning with P training examples in a two layer linear network. (b) The leading perturbative approximation to the train error breaks down when samples P becomes comparable to N . (c)-(d) Richer training reduces variance. (e)-(f) Online learning in a depth 2 linear network has identical dynamics and finite width fluctuations, but with predictor variance $\sim D/N$ for input dimension D (Appendix K).

quantify sensitivity of the kernel to perturbations in the error signal $\Delta(s)$. Lastly κ and κ_{γ} are the uncoupled variances of $K(t)$ and $K_{\gamma\gamma}(t)$ and κ_{γ} is the uncoupled covariance of $K(t)$, $K_{\gamma}(t)$.

In Fig. 3, we plot the resulting theory (diagonal blocks of q_1 from Equation 8) for two layer neural networks. As predicted by theory, all average squared deviations from the infinite width DMFT scale as $\mathcal{O}(N^{-1})$. Similarly, the average kernels $\langle K \rangle$ and test predictions $\langle f_{\gamma} \rangle$ change by a larger amount for larger γ (equation (79)). The experimental variances also match the theory quite accurately. The variance of the train error $\Delta(t)$ peaks earlier and at a lower value for richer training, but all variances go to zero at late time as the model approaches the interpolation condition $\Delta = 0$. As $\gamma \rightarrow 0$ the curve approaches $N \text{Var}(\Delta(t)) \sim \kappa y^2 t^2 e^{-2t}$, where κ is the initial NTK variance (see Section 5). While the train prediction variance goes to zero, the test point prediction does not, with richer networks reaching a lower asymptotic variance. We suspect this dynamical effect could explain lower variance observed in feature learning networks compared to lazy networks [7, 18]. In Fig. A.1, we show that the reduction in variance is not due to a reduction in the uncoupled variance $\kappa(t, s)$, which increases in γ . Rather the reduction in variance is driven by the coupling of perturbations across time.

6.2 Offline Training with Multiple Samples or Online Training in High Dimension

In this section we go beyond the single sample equations of the prior section and explore training with multiple P examples. In this case, we have training errors $\{\Delta(t)\}_{P=1}^P$ and multiple kernel entries $K(t)$ (App. E). Each of the errors $\Delta(t)$ receives a $\mathcal{O}(N^{-1/2})$ fluctuation, the training error Δ^2 has an additional variance on the order of $\mathcal{O}(\frac{P}{N})$. In the case of two-layer linear

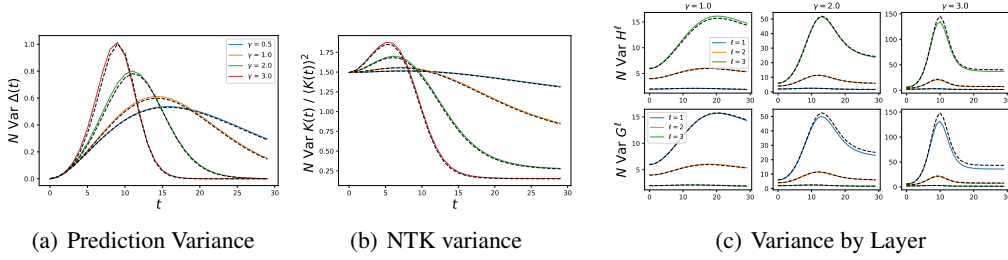


Figure 5: Depth 4 linear network with single training point. Black dashed lines are theory. (a) The variance of the training error along the task relevant subspace. We see that unlike the two layer model, more feature learning can lead to larger peaks in the finite size variance. (b) The variance of the NTK in the task relevant subspace. When properly normalized against the square of the mean $\langle K(t) \rangle^2$, the final NTK variance decreases with feature learning. (c) The gap in feature kernel variance across different layers of the network is amplified by feature learning strength γ .

networks trained on whitened data ($\frac{1}{D} \mathbf{x} \cdot \mathbf{x} = \delta$), the equations for the propagator simplify and one can separately solve for the variance of $\Delta_y(t) \in \mathbb{R}^P$ along signal direction $\mathbf{y} \in \mathbb{R}^P$ and along each of the $P - 1$ orthogonal directions (App. J). At infinite width, the task-orthogonal component Δ_\perp vanishes and only the signal dimension $\Delta_y(t)$ evolves in time with differential equation [9, 46]

$$\frac{d}{dt} \Delta_y(t) = 2 \frac{\gamma}{1 + \gamma^2 (y - \Delta_y(t))^2} \Delta_y(t), \quad \Delta_\perp(t) = 0. \quad (9)$$

However, at finite width, both the $\Delta_y(t)$ and the $P - 1$ orthogonal variables Δ_\perp inherit initialization variance, which we represent as $\Sigma_y(t, s)$ and $\Sigma_\perp(t, s)$. In Fig. 4 (a)-(b) we show this approximate solution $|\Delta(t)|^2 \sim \Delta_y(t)^2 + \frac{2}{N} \Delta_y^1(t) \Delta_y(t) + \frac{1}{N} \Sigma_y(t, t) + \frac{(P-1)}{N} \Sigma_\perp(t, t) + \mathcal{O}(N^{-2})$ across varying γ and varying P (see Appendix J for Σ_y and Σ_\perp formulas). We see that variance of train point predictions $f(t)$ increases with the total number of points despite the signal of the target vector y^2 being fixed. In this model, the bias correction $\frac{2}{N} \Delta_y^1(t) \Delta_y(t)$ is always $\mathcal{O}(1/N)$ but the variance correction is $\mathcal{O}(P/N)$. The fluctuations along the $P - 1$ orthogonal directions begin to dominate the variance at large P . Fig. 4 (b) shows that as P increases, the leading order approximation breaks down as higher order terms become relevant. Analysis for online training reveals identical fluctuation statistics, but with variance that scales as $\sim D/N$ (Appendix K) as we verify in Figure 4 (e)-(f).

7 Deep Networks

In networks deeper than two layers, the DMFT propagator has complicated dependence on non-diagonal (in time) entries of the feature kernels (see App. E). This leads to Hessian blocks with four time and four sample indices such as $D^{-1}(t, s, t^\ell, s^\ell) = \frac{\partial^4}{\partial t \partial s \partial t^\ell \partial s^\ell} \phi(h(t)) \phi(h(s))$, rendering any numerical calculation challenging. However, in deep linear networks trained on whitened data, we can exploit the symmetry in sample space and the Gaussianity of preactivation features to exactly compute derivatives without Monte Carlo sampling as we discuss in App. L. An example set of results for a depth 4 network is provided in Fig. 5. The variance for the feature kernels H^f accumulate finite size variance by layer along the forward pass and the gradient kernels G^l accumulate variance on the backward pass. The SNR of the kernels $\frac{h H_i^2}{N \text{Var}(H)}$ improves with feature learning, suggesting that richer networks will be better modeled by their mean field limits. Examples of the off-diagonal correlations obtained from the propagator are provided in App. Fig. A.3.

8 Variance can be Small Near Edge of Stability

In this section, we move beyond the gradient flow formalism and ask what large step sizes do to finite size effects. Recent studies have identified that networks trained at large learning rates can be qualitatively different than networks in the gradient flow regime, including the catapult [57] and edge of stability (EOS) phenomena [19–21]. In these settings, the kernel undergoes an initial scale

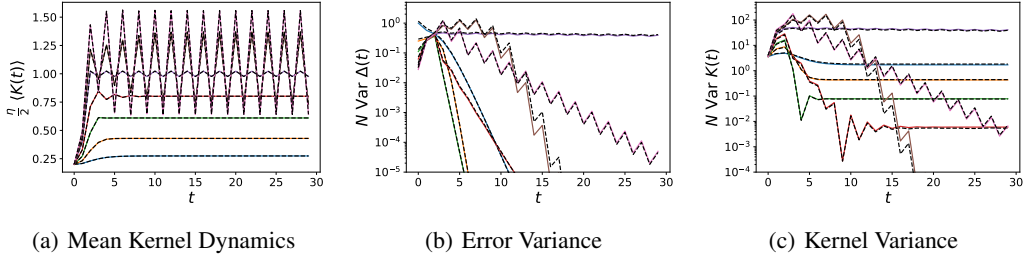


Figure 6: Edge of stability effects do not imply deviations from infinite width behavior. Black dashed lines are theory. (a) The average kernel over an ensemble of several $N = 500$ NNs (solid color). For small γ , the kernel reaches its asymptote before hitting the edge of stability. For large γ , the kernel increases and then oscillates around $2/\eta$. (b)-(c) Remarkably variance due to finite size can *reduce* during training (for γ smaller and larger than the critical value $\sim 1/\eta$), showing that infinite width DMFT can be predictive of finite NNs trained with large learning rate.

growth before exhibiting either a recovery or a clipping effect. In this section, we explore whether these dynamics are highly sensitive to initialization variance or if finite networks are well captured by mean field theory. Following [57], we consider two layer networks trained on a single example $|\mathcal{X}|^2 = D$ and $y = 1$. We use learning rate η and feature learning strength γ . The infinite width mean field equations for the prediction f_t and the kernel K_t are (App. M)

$$f_{t+1} = f_t + \eta K_t \Delta_t + \eta^2 \gamma^2 f_t \Delta_t^2, \quad K_{t+1} = K_t + 4\eta \gamma^2 f_t \Delta_t + \eta^2 \gamma^2 \Delta_t^2 K_t. \quad (10)$$

For small η , the equations are well approximated by the gradient flow limit and for small γ corresponds to a discrete time linear model. For large $\eta \gamma > 1$, the kernel K progressively sharpens (increases in scale) until it reaches $2/\eta$ and then oscillates around this value. It may be expected that near the EOS, the large oscillations in the kernels and predictions could lead to amplified finite size effects, however, we show in Fig. 6 that the leading order propagator elements decrease even after reaching the EOS threshold, indicating *reduced* disagreement between finite and infinite width dynamics.

9 Finite Width Alters Bias, Training Rate, and Variance in Realistic Tasks

To analyze the effect of finite width on neural network dynamics during realistic learning tasks, we studied a vanilla depth-6 ReLU CNN trained on CIFAR-10 (experimental details in App. B, G.2) In Fig. 7, we train an ensemble of $E = 8$ independently initialized CNNs of each width N . Wider networks not only have better performance for a single model (solid), but also have lower bias (dashed), measured with ensemble averaging of the logits. Because of faster convergence of wide networks, we observe wider networks have higher variance, but if we plot variance at fixed ensemble training accuracy, wider networks have consistently lower variance (Fig. 7(d)).

We next seek an explanation for why wider networks after ensembling trains at a faster *rate*. Theoretically, this can be rationalized by a finite-width alteration to the ensemble averaged NTK, which governs the convergence timescale of the ensembled predictions (App. G.1). Our analysis in App. G.1 suggests that the rate of convergence receives a finite size correction with leading correction $\mathcal{O}(N^{-1})$ G.2. To test this hypothesis, we fit the ensemble training loss curve to exponential function $\mathcal{L} \approx C \exp(-R_N t)$ where C is a constant. We plot the fit R_N as a function of N^{-1} result in Fig. 7(e). For large N , we see the leading behavior is linear in N^{-1} , but begins to deviate at small N as a quadratic function of N^{-1} , suggesting that second order effects become relevant around $N \lesssim 100$.

In App. Fig. A.4, we train a smaller subset of CIFAR-10 where we find that R_N is well approximated by a $\mathcal{O}(N^{-1})$ correction, consistent with the idea that higher sample size drives the dynamics out of the leading order picture. We also analyze the effect of γ on variance in this task. In App. Fig. A.5, we train $N = 64$ models with varying γ . Increased γ reduces variance of the logits and alters the representation (measured with kernel-task alignment), the training and test accuracy are roughly insensitive to the richness γ in the range we considered.

10 Discussion

We studied the leading order fluctuations of kernels and predictions in mean field neural networks. Feature learning dynamics can reduce undesirable finite size variance, making finite networks order

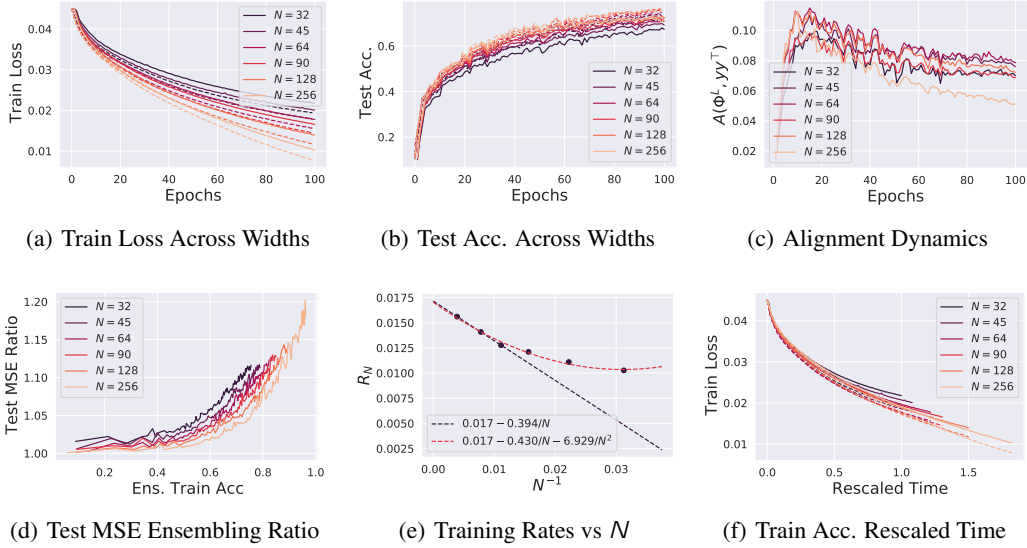


Figure 7: Depth 6 CNN trained on CIFAR-10 for different widths N with richness $\gamma = 0.2$, $E = 8$ ensembles. (a)-(b) For this range of widths, we find that smaller networks perform worse in train and test error, not only in terms of the single models (solid) but also in terms of bias (dashed). The delayed training of ensembled finite width models indicates that the correction to the mean order parameters (App. G) is non-negligible. (c) Alignment of the average kernel to test labels is also not conserved across width. (d) The ratio of the test MSE for a single model to the ensembled logit MSE. (e) The fitted rate R_N of training width N models as a function of N^{-1} . We rescale the time axis by R_N to allow for a fair comparison of prediction variance for networks at comparable performance levels. (f) In rescaled time, ensembled network training losses (dashed) are coincident.

parameters closer to the infinite width limit. In several toy models, we revealed some interesting connections between the influence of feature learning, depth, sample size, and large learning rate and the variance of various DMFT order parameters. Lastly, in realistic tasks, we illustrated that bias corrections can be significant as rates of learning can be modified by width. Though our full set of equations for the leading finite size fluctuations are quite general in terms of network architecture and data structure, they are only derived at the level of rigor of physics rather than a formally rigorous proof which would need several additional assumptions to make the perturbation expansion properly defined. Further, the leading terms in our perturbation series involving only μ does not capture the complete finite size distribution defined in Eq. (3), especially as the sample size becomes comparable to the width. It would be interesting to see if proportional limits of the rich training regime where samples and width scale linearly can be examined dynamically [58]. Future work could explore in greater detail the higher order contributions from averages involving powers of $V(q)$ by examining cubic and higher derivatives of S in Eq. (3). It could also be worth examining in future work how finite size impacts other biologically plausible learning rules, where the effective NTK can have asymmetric (over sample index) fluctuations [46]. Also of interest would be computing the finite width effects in other types of architectures, including residual networks with various branch scalings [59, 60]. Further, even though we expect our perturbative expressions to give a precise asymptotic description of finite networks in mean field/ μ P, the resulting expressions are not realistically computable in deep networks trained on large dataset size P for long times T since the number of Hessian entries scales as $\mathcal{O}(T^4 P^4)$ and a matrix of this size must be stored in memory and inverted in the general case. Future work could explore solvable special cases such as high dimensional limits.

Code Availability

Code to reproduce the experiments in this paper is provided at https://github.com/Pehlevan-Group/dmft_fluctuations. Details about numerical methods and computational implementation can be found in Appendices F and N.

Acknowledgements

CP is supported by NSF Award DMS-2134157, NSF CAREER Award IIS-2239780, and a Sloan Research Fellowship. BB is supported by a Google PhD research fellowship and NSF Award DMS-2134157. This work has been made possible in part by a gift from the Chan Zuckerberg Initiative Foundation to establish the Kempner Institute for the Study of Natural and Artificial Intelligence. The computations in this paper were run on the FASRC cluster supported by the FAS Division of Science Research Computing Group at Harvard University. BB thanks Alex Atanasov, Jacob Zavatone-Veth for their comments on this manuscript and Boris Hanin, Greg Yang, Mufan Bill Li and Jeremy Cohen for helpful discussions.

References

- [1] Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31, pages 8571–8580. Curran Associates, Inc., 2018.
- [2] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. *Advances in Neural Information Processing Systems*, 32, 2019.
- [3] Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. *Advances in neural information processing systems*, 32, 2019.
- [4] Lenaic Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. *Advances in Neural Information Processing Systems*, 32, 2019.
- [5] Greg Yang and Etai Littwin. Tensor programs iib: Architectural universality of neural tangent kernel training dynamics. In *International Conference on Machine Learning*, pages 11762–11772. PMLR, 2021.
- [6] Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit. In *Conference on Learning Theory*, pages 2388–2464. PMLR, 2019.
- [7] Mario Geiger, Stefano Spigler, Arthur Jacot, and Matthieu Wyart. Disentangling feature and lazy training in deep neural networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2020(11):113301, 2020.
- [8] Greg Yang and Edward J Hu. Tensor programs iv: Feature learning in infinite-width neural networks. In *International Conference on Machine Learning*, pages 11727–11737. PMLR, 2021.
- [9] Blake Bordelon and Cengiz Pehlevan. Self-consistent dynamical field theory of kernel evolution in wide neural networks. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [10] Jiaoyang Huang and Horng-Tzer Yau. Dynamics of deep neural networks and neural tangent hierarchy. In *International conference on machine learning*, pages 4542–4551. PMLR, 2020.
- [11] Ethan Dyer and Guy Gur-Ari. Asymptotics of wide networks from feynman diagrams. In *International Conference on Learning Representations*, 2020.
- [12] Anders Andreassen and Ethan Dyer. Asymptotics of wide convolutional neural networks. *arXiv preprint arXiv:2008.08675*, 2020.
- [13] Daniel A Roberts, Sho Yaida, and Boris Hanin. *The principles of deep learning theory*. Cambridge University Press Cambridge, MA, USA, 2022.
- [14] Boris Hanin. Random fully connected neural networks as perturbatively solvable hierarchies. *arXiv preprint arXiv:2204.01058*, 2022.

- [15] Sho Yaida. Meta-principled family of hyperparameter scaling strategies. *arXiv preprint arXiv:2210.04909*, 2022.
- [16] Jaehoon Lee, Samuel Schoenholz, Jeffrey Pennington, Ben Adlam, Lechao Xiao, Roman Novak, and Jascha Sohl-Dickstein. Finite versus infinite neural networks: an empirical study. *Advances in Neural Information Processing Systems*, 33:15156–15172, 2020.
- [17] Greg Yang, Edward Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ryder, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tuning large neural networks via zero-shot hyperparameter transfer. *Advances in Neural Information Processing Systems*, 34, 2021.
- [18] Alexander Atanasov, Blake Bordelon, Sabarish Sainathan, and Cengiz Pehlevan. The onset of variance-limited behavior for networks in the lazy and rich regimes. In *The Eleventh International Conference on Learning Representations*, 2023.
- [19] Jeremy Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent on neural networks typically occurs at the edge of stability. In *International Conference on Learning Representations*, 2021.
- [20] Alex Damian, Eshaan Nichani, and Jason D. Lee. Self-stabilization: The implicit bias of gradient descent at the edge of stability. In *The Eleventh International Conference on Learning Representations*, 2023.
- [21] Atish Agarwala, Fabian Pedregosa, and Jeffrey Pennington. Second-order regression models exhibit progressive sharpening to the edge of stability. *arXiv preprint arXiv:2210.04860*, 2022.
- [22] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research).
- [23] Radford M Neal. Priors for infinite networks. In *Bayesian Learning for Neural Networks*, pages 29–53. Springer, 1996.
- [24] Radford M Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media, 2012.
- [25] Jaehoon Lee, Jascha Sohl-dickstein, Jeffrey Pennington, Roman Novak, Sam Schoenholz, and Yasaman Bahri. Deep neural networks as gaussian processes. In *International Conference on Learning Representations*, 2018.
- [26] Alexander G. de G. Matthews, Jiri Hron, Mark Rowland, Richard E. Turner, and Zoubin Ghahramani. Gaussian process behaviour in wide deep neural networks. In *International Conference on Learning Representations*, 2018.
- [27] Boris Hanin and Mihai Nica. Finite depth and width corrections to the neural tangent kernel. In *International Conference on Learning Representations*, 2020.
- [28] Boris Hanin and Mihai Nica. Products of many large random matrices and gradients in deep neural networks. *Communications in Mathematical Physics*, 376(1):287–322, 2020.
- [29] Sho Yaida. Non-gaussian processes and neural networks at finite widths. In *Mathematical and Scientific Machine Learning*, pages 165–192. PMLR, 2020.
- [30] Jacob Zavatone-Veth and Cengiz Pehlevan. Exact marginal prior distributions of finite bayesian neural networks. *Advances in Neural Information Processing Systems*, 34:3364–3375, 2021.
- [31] Jacob Zavatone-Veth, Abdulkadir Canatar, Ben Ruben, and Cengiz Pehlevan. Asymptotics of representation learning in finite bayesian neural networks. *Advances in Neural Information Processing Systems*, 34, 2021.
- [32] Gadi Naveh, Oded Ben David, Haim Sompolinsky, and Zohar Ringel. Predicting the outputs of finite deep neural networks trained with noisy gradients. *Physical Review E*, 104(6):064301, 2021.

- [33] Adam X. Yang, Maxime Robeyns, Edward Milsom, Ben Anson, Nandi Schoots, and Laurence Aitchison. A theory of representation learning gives a deep generalisation of kernel methods. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 39380–39415. PMLR, 23–29 Jul 2023.
- [34] Inbar Seroussi, Gadi Naveh, and Zohar Ringel. Separation of scales and a thermodynamic description of feature learning in some cnns. *Nature Communications*, 14(1):908, 2023.
- [35] Qianyi Li and Haim Sompolinsky. Statistical mechanics of deep linear neural networks: The backpropagating kernel renormalization. *Physical Review X*, 11(3):031059, 2021.
- [36] Jacob A Zavatore-Veth and Cengiz Pehlevan. Depth induces scale-averaging in overparameterized linear bayesian neural networks. *55th Asilomar Conference on Signals, Systems, and Computers*, 2021.
- [37] Jacob A Zavatore-Veth, William L Tong, and Cengiz Pehlevan. Contrasting random and learned features in deep bayesian linear regression. *Physical Review E*, 105(6):064118, 2022.
- [38] Gadi Naveh and Zohar Ringel. A self consistent theory of gaussian processes captures feature learning effects in finite cnns. *Advances in Neural Information Processing Systems*, 34, 2021.
- [39] Jacob A Zavatore-Veth, Abdulkadir Canatar, Benjamin S Ruben, and Cengiz Pehlevan. Asymptotics of representation learning in finite bayesian neural networks*. *Journal of Statistical Mechanics: Theory and Experiment*, 2022(11):114008, nov 2022.
- [40] Grant Rotskoff and Eric Vanden-Eijnden. Parameters as interacting particles: long time convergence and asymptotic error scaling of neural networks. *Advances in neural information processing systems*, 31, 2018.
- [41] Lenaic Chizat and Francis Bach. On the global convergence of gradient descent for overparameterized models using optimal transport. *Advances in neural information processing systems*, 31, 2018.
- [42] Dyego Araújo, Roberto I Oliveira, and Daniel Yukimura. A mean-field limit for certain deep neural networks. *arXiv preprint arXiv:1906.00193*, 2019.
- [43] Rodrigo Veiga, Ludovic Stephan, Bruno Loureiro, Florent Krzakala, and Lenka Zdeborova. Phase diagram of stochastic gradient descent in high-dimensional two-layer neural networks. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [44] Luca Arnaboldi, Ludovic Stephan, Florent Krzakala, and Bruno Loureiro. From high-dimensional amp; mean-field dynamics to dimensionless odes: A unifying approach to sgd in two-layers networks. In Gergely Neu and Lorenzo Rosasco, editors, *Proceedings of Thirty Sixth Conference on Learning Theory*, volume 195 of *Proceedings of Machine Learning Research*, pages 1199–1227. PMLR, 12–15 Jul 2023.
- [45] Huy Tuan Pham and Phan-Minh Nguyen. Limiting fluctuation and trajectorial stability of multilayer neural networks with mean field training. *Advances in Neural Information Processing Systems*, 34:4843–4855, 2021.
- [46] Blake Bordelon and Cengiz Pehlevan. The influence of learning rule on representation dynamics in wide neural networks. In *The Eleventh International Conference on Learning Representations*, 2023.
- [47] Paul Cecil Martin, ED Siggia, and HA Rose. Statistical dynamics of classical systems. *Physical Review A*, 8(1):423, 1973.
- [48] Moshe Moshe and Jean Zinn-Justin. Quantum field theory in the large n limit: A review. *Physics Reports*, 385(3-6):69–228, 2003.

- [49] Jean Zinn-Justin. *Quantum field theory and critical phenomena*, volume 171. Oxford university press, 2021.
- [50] Carson C Chow and Michael A Buice. Path integral methods for stochastic differential equations. *The Journal of Mathematical Neuroscience (JMN)*, 5:1–35, 2015.
- [51] Moritz Helias and David Dahmen. *Statistical Field Theory for Neural Networks*. Springer International Publishing, 2020.
- [52] A Crisanti and H Sompolinsky. Path integral approach to random neural networks. *Physical Review E*, 98(6):062120, 2018.
- [53] Kai Segadlo, Bastian Epping, Alexander van Meegen, David Dahmen, Michael Krämer, and Moritz Helias. Unified field theoretical approach to deep and recurrent neuronal networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2022(10):103401, 2022.
- [54] Yizhang Lou, Chris E Mingard, and Soufiane Hayou. Feature learning and signal propagation in deep neural networks. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 14248–14282. PMLR, 17–23 Jul 2022.
- [55] Carl M Bender and Steven Orszag. *Advanced mathematical methods for scientists and engineers I: Asymptotic methods and perturbation theory*, volume 1. Springer Science & Business Media, 1999.
- [56] Mehran Kardar. *Statistical physics of fields*. Cambridge University Press, 2007.
- [57] Aitor Lewkowycz, Yasaman Bahri, Ethan Dyer, Jascha Sohl-Dickstein, and Guy Gur-Ari. The large learning rate phase of deep learning: the catapult mechanism. *arXiv preprint arXiv:2003.02218*, 2020.
- [58] Ben Adlam and Jeffrey Pennington. The neural tangent kernel in high dimensions: Triple descent and a multi-scale theory of generalization. In *International Conference on Machine Learning*, pages 74–84. PMLR, 2020.
- [59] Soufiane Hayou. On the infinite-depth limit of finite-width neural networks. *Transactions on Machine Learning Research*, 2023.
- [60] Blake Bordelon, Lorenzo Noci, Mufan Bill Li, Boris Hanin, and Cengiz Pehlevan. Depthwise hyperparameter transfer in residual networks: Dynamics and scaling limit. *arXiv preprint arXiv:2309.16620*, 2023.
- [61] Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas Steiner, and Marc van Zee. Flax: A neural network library and ecosystem for JAX, 2023.
- [62] Michael Baake and Ulrike Schlaegel. The peano-baker series. *Proceedings of the Steklov Institute of Mathematics*, 275(1):155–159, 2011.
- [63] Roger W Brockett. *Finite dimensional linear systems*. SIAM, 2015.
- [64] Alexander Atanasov, Blake Bordelon, and Cengiz Pehlevan. Neural networks as kernel learners: The silent alignment effect. In *International Conference on Learning Representations*, 2022.
- [65] B. Bordelon, A. Canatar, and C. Pehlevan. Spectrum dependent learning curves in kernel regression and wide neural networks. *International Conference of Machine Learning*, 2020.
- [66] Phan-Minh Nguyen. Mean field limit of the learning dynamics of multilayer neural networks. *arXiv preprint arXiv:1902.02880*, 2019.

Appendix

A Additional Figures

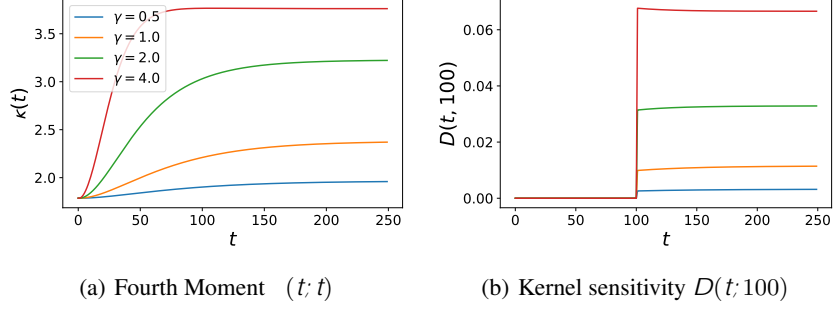


Figure A.1: The κ and D functions for varying γ in Figure 3. (a) The uncoupled kernel variance $\kappa(t, t)$ increases monotonically with γ . This reveals that the dynamical filtering of κ is what is responsible for the late time decrease in variance during feature learning. (b) The tensor $D(t, s)$ measures sensitivity of kernel at time t to perturbation in Δ at time s . The $D(t, s)$ entries also increase with γ . This suggests that the reduction in variance of the training error and the kernel are not due to reduction in κ , but rather a dynamical filtering effect due to scale growth in K_γ and rapid reduction in Δ_γ .

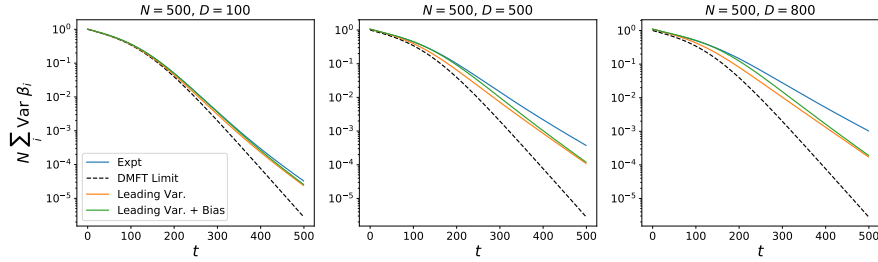


Figure A.2: A comparison of the bias and variance corrections in the toy model of Figure 4. At small D/N (or P/N for offline training) the leading variance and the variance and bias both track the experiment. Both the bias and the variance contribute positively towards the total generalization error since the variance correction alone (orange) exceeds the DMFT limiting error (dashed) and the variance and bias correction together (green) exceed variance alone (orange). However, for large D/N (or P/N) the leading order picture fails to describe the finite width experiment, indicating significant variance possibly at higher order scales (like $D^2/N^2, D^3/N^3, \dots$).

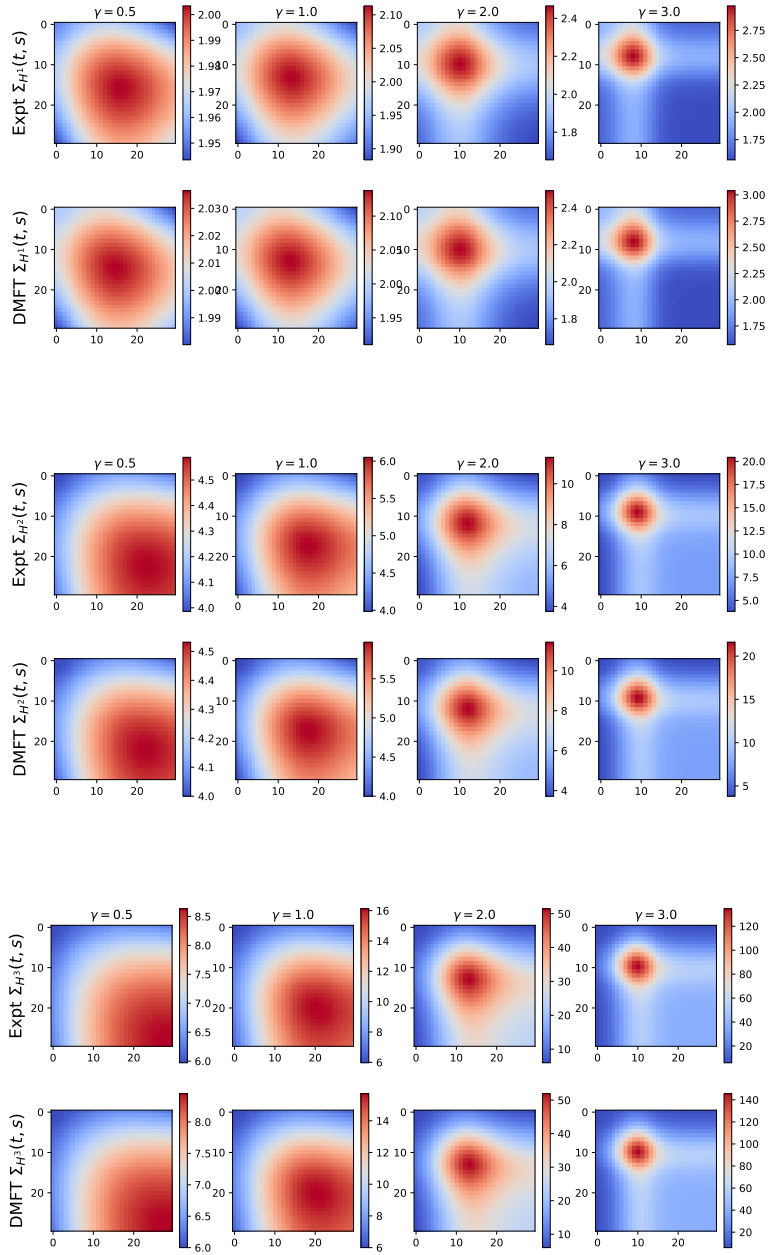


Figure A.3: The covariance of kernel entries across pairs of time points $\Sigma_{H^i}(t, s) = N \text{Cov}(H^i(t, t), H^i(s, s))$ for depth 4 linear network trained on whitened data. The variance becomes increasingly localized in time as feature learning γ increases.

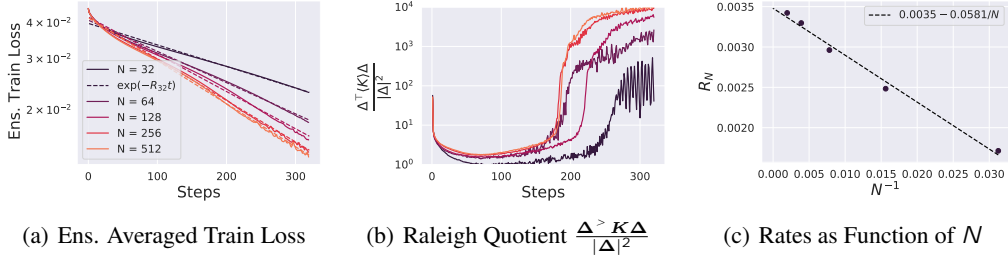


Figure A.4: The ensemble averaged train loss for the same depth 6 CNN trained on a random subsample of $P = 64$ CIFAR-10 points. Training is full batch gradient descent with $\gamma = 0.05$. (a) The ensemble train accuracy for this subset of CIFAR-10 is well modeled as an exponential in time $\mathcal{L}(t) \propto \exp(-R_N t)$ with a rate R_N that depends on width. (b) The projection of the errors Δ on the average NTK $\langle K \rangle$ (which is related to the rate of decay of the training loss, see Appendix G) reveals that wider networks are more aligned with their instantaneous error signals. (c) The rates R_N are indeed a linear functions of N^{-1} , with $R_N = R_\gamma + \frac{R^1}{N}$, consistent with the average NTK $\langle K \rangle$ receiving a N^{-1} correction. Using ensembling to find a scaling law like that above can thus allow one extrapolate the training rate of infinite width mean field models.

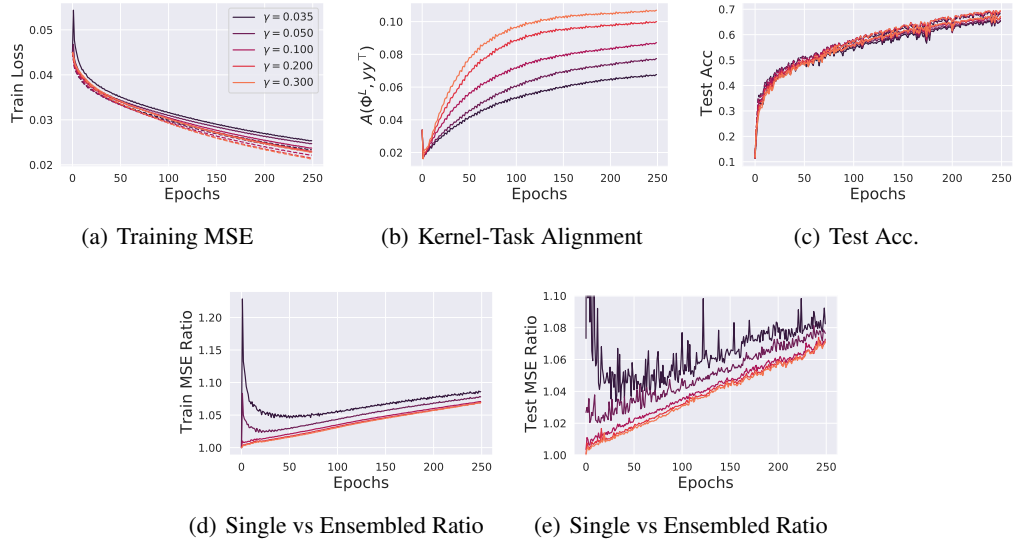


Figure A.5: Width $N = 64$ depth 6 CNNs trained on the full CIFAR-10 with MSE. An ensemble of size $E = 10$ randomly initialized networks are trained. (a) Training MSE for varying γ . (b) Final layer kernel-task alignment does strongly depend on γ , despite similar train dynamics. (c) Top-1 classification test accuracy is only slightly different across γ . A small benefit from ensembling is visible late in training. (d) Initialization variance (measured by the ratio of single model to ensembled MSE) for training and test losses. Richer networks have lower variance throughout training. (e) Networks have distinct kernel dynamics when trained with different γ as evidenced by the alignment (cosine similarity) between the final layer feature kernel Φ^L and the target test labels \mathbf{y} .

B CIFAR-10 Experimental Details

We trained the following depth 6 CNN architecture in the mean field parameterization using FLAX [61] on a single GPU. The bias parameters were zero in each hidden Conv layer, but were used for the readout weights. The networks were trained with MSE loss on centered 10 dimensional targets $\mathbf{y} \in \mathbb{R}^{10}$ for $\mu \in [P]$. Each convolution was followed by an average pooling operation. To obtain mean field behavior, NTK parameterization with a modified final layer is used [7, 9].

```
1 from flax import linen as nn
2 import jax.numpy as jnp
3
4 class CNN(nn.Module):
5
6     width: int
7
8     def setup(self):
9         kif = nn.initializers.normal(stddev = 1.0) # 0_N(1) entries
10        self.conv1 = nn.Conv(features = self.width, kernel_init = kif,
11        use_bias = False, kernel_size = (3,3))
12        self.conv2 = nn.Conv(features = self.width, kernel_init = kif,
13        use_bias = False, kernel_size = (3,3))
14        self.conv3 = nn.Conv(features = self.width, kernel_init = kif,
15        use_bias = False, kernel_size = (3,3))
16        self.conv4 = nn.Conv(features = self.width, kernel_init = kif,
17        use_bias = False, kernel_size = (3,3))
18        self.conv5 = nn.Conv(features = self.width, kernel_init = kif,
19        use_bias = False, kernel_size = (3,3))
20        self.readout = nn.Dense(features = 10, use_bias = True,
21        kernel_init = kif)
22        return
23
24    def __call__(self, x, train = True):
25        N = self.width
26        D = 3
27        x = self.conv1(x) / jnp.sqrt(D * 9)
28        x = jnp.sqrt(2.0) * nn.relu(x)
29        x = nn.avg_pool(x, window_shape=(2,2), strides = (2,2)) # 32 x 32
30        -> 16 x 16
31        x = self.conv2(x) / jnp.sqrt(N*9) # explicit N^{-1/2}
32        x = jnp.sqrt(2.0) * nn.relu(x)
33        x = nn.avg_pool(x, window_shape=(2,2), strides = (2,2)) # 16 x 16
34        -> 8 x 8
35        x = self.conv3(x)/jnp.sqrt(N*9)
36        x = jnp.sqrt(2.0) * nn.relu(x)
37        x = nn.avg_pool(x, window_shape=(2,2), strides = (2,2)) # 8 x 8 ->
38        4 x 4
39        x = self.conv4(x) / jnp.sqrt(N*9)
40        x = jnp.sqrt(2.0) * nn.relu(x)
41        x = nn.avg_pool(x, window_shape=(2,2), strides = (2,2)) # 4 x 4
42        -> 2 x 2
43        x = self.conv5(x) / jnp.sqrt(N*9)
44        x = jnp.sqrt(2.0) * nn.relu(x)
45        x = nn.avg_pool(x, window_shape=(2,2), strides = (2,2)) # 2 x 2
46        -> 1 x 1
47        x = x.reshape((x.shape[0], -1)) # flatten
48        x = self.readout(x) / N # for mean field scaling
49        return x
```

All models were trained with standard SGD with a batch size of 256. Each element in the ensemble of E networks is trained on identical batches presented in identical order. For the Figure 7 experiments, the raw learning rate is scaled as $\eta = 10N\sqrt{\gamma}$ with $\gamma = 0.2$ (note that mean field theory requires scaling the raw learning rate linearly with N since the raw NTK is $\mathcal{O}(N^{-1})$ [9]). For Figure A.5, the learning rate is $\eta = 5N\sqrt{\gamma}$. We find that choosing $\eta \propto \sqrt{\gamma}$ gives approximately conserved training

times across γ (though distinct representation dynamics). The Figure A.4 shows the dynamics of fitting $P = 64$ training points with full batch gradient descent and $\gamma = 0.1$.

C Review of DMFT: Deriving the Action

In this section we derive the DMFT action which contains all of the necessary statistical information about randomly initialized finite width N networks. From the action S the DMFT saddle point and the propagator can be computed. This derivation follows closely the original derivation by Bordelon & Pehlevan [9]. We start by writing the gradient flow dynamics on weight matrices

$$\frac{d}{dt} \mathbf{W}(t) = \frac{\gamma}{\sqrt{N}} \sum_{=1}^{\mathcal{X}} \Delta(t) \mathbf{g}^{+1}(t) \phi(\mathbf{h}(t)) \quad (11)$$

where $\Delta(t) = -\frac{\partial L}{\partial \mathbf{f}(t)}$ are the error signals and $\mathbf{g}(t) = N\gamma \frac{\partial \mathbf{f}}{\partial \mathbf{h}(t)}$ are the back-propagation signals. The prediction dynamics satisfy

$$\frac{d}{dt} \mathbf{f}(t) = \sum^{\times} K(t) \Delta(t) \quad (12)$$

where $K(t)$ is the instantaneous neural tangent kernel (NTK). At finite width N all of the above quantities depend on the precise initialization of the network. We transform the weight dynamics into an integral equation and use the recurrence for \mathbf{h} to obtain the following

$$\begin{aligned} \mathbf{h}^{+1}(t) &= \frac{1}{\sqrt{N}} \mathbf{W}(0) \phi(\mathbf{h}(t)) + \gamma \int_0^t ds \sum^{\times} \Phi(t, s) \mathbf{g}^{+1}(s) \\ \mathbf{g}(t) &= \dot{\phi}(\mathbf{h}(t)) \odot \mathbf{z}(t) \\ \mathbf{z}(t) &= \frac{1}{\sqrt{N}} \mathbf{W}(0) \mathbf{g}^{+1}(t) + \gamma \int_0^t ds \sum^{\times} G^{+1}(t, s) \phi(\mathbf{h}(s)) \end{aligned} \quad (13)$$

where we introduced the feature and gradient kernels

$$\Phi(t, s) = \frac{1}{N} \phi(\mathbf{h}(t)) \cdot \phi(\mathbf{h}(s)), \quad G(t, s) = \frac{1}{N} \mathbf{g}(t) \cdot \mathbf{g}(s). \quad (14)$$

Written this way, we see that the source of the disorder which depends on the initial random weights $\mathbf{W}(0)$ comes through the fields

$$\mathbf{h}^{+1}(t) = \frac{1}{\sqrt{N}} \mathbf{W}(0) \phi(\mathbf{h}(t)), \quad \mathbf{z}^{+1}(t) = \frac{1}{\sqrt{N}} \mathbf{W}(0) \mathbf{g}(t). \quad (15)$$

If we can characterize the distribution of the fields $\mathbf{h}(t)$ and $\mathbf{z}(t)$, then we can consequently characterize the distribution of $\mathbf{h}(t), \mathbf{g}(t)$. We therefore choose to study the moment generating functional

$$Z[\{\mathbf{j}, \mathbf{v}\}] = \exp \left[\int dt \sum^{\times} \mathbf{j}(t) \cdot \mathbf{h}(t) + \mathbf{v}(t) \cdot \mathbf{z}(t) \right] \quad (16)$$

Moments of these fields can be computed through differentiation with respect to the sources \mathbf{j}, \mathbf{v} near zero-source ($\mathbf{j} = \mathbf{v} = 0$)

$$\begin{aligned} &\chi_1^1(t_1) \dots \chi_n^n(t_n) \xi_1^1(t_1) \dots \xi_m^m(t_m) \\ &= \frac{\delta}{\delta j_1^1(t_1)} \dots \frac{\delta}{\delta j_n^n(t_n)} \frac{\delta}{\delta v_1^1(t_1)} \dots \frac{\delta}{\delta v_m^m(t_m)} Z[\{\mathbf{j}, \mathbf{v}\}]|_{\mathbf{j}=\mathbf{v}=0}. \end{aligned} \quad (17)$$

To average over the initial weights, we introduce a Fourier representation of the Dirac-Delta function $1 = \int dz \delta(z) = \int \frac{dz d\bar{z}}{2} \exp(i\bar{z}z)$. We perform this transformation for each of the fields to enforce

their definition

$$\begin{aligned}\delta \left(\hat{h}(t) - \frac{1}{\sqrt{N}} \mathbf{W}(0) \phi(\mathbf{h}(t)) \right) &= \int \frac{d\hat{h}(t)}{(2\pi)^N} \exp \left(i \hat{h}(t) \cdot \left(\hat{h}(t) - \frac{1}{\sqrt{N}} \mathbf{W}(0) \phi(\mathbf{h}(t)) \right) \right) \\ \delta \left(\hat{h}(t) - \frac{1}{\sqrt{N}} \mathbf{W}(0) \mathbf{g}^{i+1}(t) \right) &= \int \frac{d\hat{h}(t)}{(2\pi)^N} \exp \left(i \hat{h}(t) \cdot \left(\hat{h}(t) - \frac{1}{\sqrt{N}} \mathbf{W}(0) \mathbf{g}^{i+1}(t) \right) \right)\end{aligned}\quad (18)$$

We insert these Dirac delta functions so that we can directly average over the weights

$$\begin{aligned}\ln E_{\mathbf{W}(0)} \exp \left[-\frac{i}{\sqrt{N}} \text{Tr} \mathbf{W}(0) \int dt \int \mathbf{h}^{i+1}(t) \phi(\mathbf{h}(t)) \cdot \mathbf{g}^{i+1}(t) \right] \\ = -\frac{1}{2} \int dt ds \int \mathbf{h}^{i+1}(t) \cdot \mathbf{h}^{i+1}(s) \Phi(t, s) + \int dt ds \hat{h}^{i+1}(t) \cdot \hat{h}^{i+1}(s) G^{i+1}(t, s) \\ - \frac{1}{N} \int dt ds \left(\hat{h}^{i+1}(t) \cdot \mathbf{g}^{i+1}(s) \right) \left(\phi(\mathbf{h}(t)) \cdot \hat{h}(s) \right)\end{aligned}\quad (19)$$

where we introduced the kernels Φ, G . We next introduce the order parameter

$$A(t, s) = -\frac{i}{N} \phi(\mathbf{h}(t)) \cdot \hat{h}(s) \quad (20)$$

To enforce the definitions of the new order parameters $\{\Phi, G, A\}$ we again introduce Dirac-delta functions

$$\begin{aligned}\delta \left(N \Phi(t, s) - \phi(\mathbf{h}(t)) \cdot \phi(\mathbf{h}(s)) \right) \\ = \int \frac{d\hat{\Phi}(t, s)}{2\pi i} \exp \left(i \hat{\Phi}(t, s) \left(N \Phi(t, s) - \phi(\mathbf{h}(t)) \cdot \phi(\mathbf{h}(s)) \right) \right)\end{aligned}\quad (21)$$

Analogous constraints for G and A are enforced with conjugate variables \hat{G}, B . After introducing these variables, we find that the moment generating functional has the form

$$\begin{aligned}Z = \int \prod_{ts} \frac{d\hat{\Phi}(t, s) d\Phi(t, s) d\hat{G}(t, s) dG(t, s) d\hat{G}(t, s) dG(t, s) dA(t, s) dB(t, s)}{2\pi i} \\ \exp \left[NS \left[\{\Phi, \hat{\Phi}, G, \hat{G}, A, B\} \right] \right]\end{aligned}\quad (22)$$

where S is the $\mathcal{O}(1)$ DMFT *action* which defines the statistical distribution over the dynamics. The action takes the form

$$\begin{aligned}S = \int dt ds \int \mathbf{h} \left(\hat{\Phi}(t, s) \Phi(t, s) + \hat{G}(t, s) - A(t, s) B(s, t) \right) \\ + \frac{1}{N} \sum_{i=1}^N \ln \mathcal{Z}[\{j_i, v_i\}]\end{aligned}\quad (23)$$

where \mathcal{Z} is the single site stochastic process for layer ℓ which defines the marginal distribution of χ, ξ , with the following form

$$\begin{aligned}\mathcal{Z}[\{j(t), v(t)\}] = \int_t \frac{d\hat{\chi}(t) d\chi(t) d\hat{\xi}(t) d\xi(t)}{2\pi} \exp \left[\int dt \int \mathbf{h} \left[j(t) \chi(t) + v(t) \xi(t) \right] \right] \\ \exp \left[-\frac{1}{2} \int dt ds \int \mathbf{h} \left(\hat{\Phi}(t, s) \hat{\chi}(t) \hat{\chi}(s) + \hat{G}(t, s) \hat{\xi}(t) \hat{\xi}(s) \right) \right] \\ \exp \left[-i \int dt ds \int \mathbf{h} \left(B(t, s) \hat{\xi}(t) \phi(\mathbf{h}(s)) + A^{-1}(t, s) \hat{\chi}(t) g(s) \right) \right] \\ \exp \left[i \int dt \left[\hat{\chi}(t) \chi(t) + \hat{\xi}(t) \xi(t) \right] \right]\end{aligned}\quad (24)$$

where in the above, the $\{h, g\}$ fields should be regarded as functionals of $\{\chi, \xi\}$. At zero source $\mathbf{j}, \mathbf{v} \rightarrow 0$ this function S can be regarded as the log density for the complete collection of order parameters $\mathbf{q} = \{\hat{\Phi}, \Phi, \hat{G}, G, A, B\}$ which collectively control the dynamics. Concretely, we have that $p(\mathbf{q}) \propto \exp(NS(\mathbf{q}))$. In the next section we explore an approximation scheme for averages over this distribution at large N .

D Cumulant Expansion of Observables

We are interested in a principled power series expansion (in $1/N$) of any observable average $\langle O(\mathbf{q}) \rangle$ that depends on DMFT order parameters \mathbf{q} . At any width N the observable average takes the form

$$\langle O(\mathbf{q}) \rangle_N = \frac{\int_{\mathbf{q}} d\mathbf{q} \exp(NS(\mathbf{q})) O(\mathbf{q})}{\int_{\mathbf{q}} d\mathbf{q} \exp(NS(\mathbf{q}))} \quad (25)$$

As discussed in the main text, the $N \rightarrow \infty$ limit gives $\langle O(\mathbf{q}) \rangle_N \sim O(\mathbf{q}_1)$ where $\left. \frac{\partial S}{\partial \mathbf{q}} \right|_{\mathbf{q}_1} = 0$ by a steepest descent argument [55]. We assume that S 's Hessian is negative semidefinite so that $\left. -\nabla^2 S(\mathbf{q}) \right|_{\mathbf{q}_1} \succeq 0$ and Taylor expand $S(\mathbf{q})$ around the saddle point \mathbf{q}_1 giving $S(\mathbf{q}) = S(\mathbf{q}_1) + \frac{1}{2}(\mathbf{q} - \mathbf{q}_1)^T \nabla^2 S(\mathbf{q})(\mathbf{q} - \mathbf{q}_1) + V(\mathbf{q} - \mathbf{q}_1)$. We note that the remainder function V contains only cubic and higher powers of $\mathbf{q} - \mathbf{q}_1 \equiv \delta \mathbf{q} / \sqrt{N}$. The variable $\delta \mathbf{q}$ will be order $\mathcal{O}(1)$. This will allow us to verify that additional terms are suppressed in powers of $1/N$. Expanding both the numerator and denominator's integrands in powers of V , we find

$$\begin{aligned} \langle O(\mathbf{q}) \rangle_N &= \frac{\int_{\mathbf{q}} d\mathbf{q} \exp\left[-\frac{N}{2}(\mathbf{q} - \mathbf{q}_1)^T \nabla^2 S(\mathbf{q})(\mathbf{q} - \mathbf{q}_1) + NV(\mathbf{q} - \mathbf{q}_1)\right] O(\mathbf{q})}{\int_{\mathbf{q}} d\mathbf{q} \exp\left[-\frac{N}{2}(\mathbf{q} - \mathbf{q}_1)^T \nabla^2 S(\mathbf{q})(\mathbf{q} - \mathbf{q}_1) + NV(\mathbf{q} - \mathbf{q}_1)\right]} \\ &= \frac{\int_{\mathbf{q}} d\mathbf{q} \exp\left[-\frac{1}{2} \delta \mathbf{q}^T \nabla^2 S(\mathbf{q}_1) \delta \mathbf{q} + N \delta \mathbf{q}^T V(\mathbf{q}_1) + \dots\right] O(\mathbf{q}_1 + \delta \mathbf{q})}{\int_{\mathbf{q}} d\mathbf{q} \exp\left[-\frac{1}{2} \delta \mathbf{q}^T \nabla^2 S(\mathbf{q}_1) \delta \mathbf{q} + N \delta \mathbf{q}^T V(\mathbf{q}_1) + \dots\right]} \\ &= \frac{\langle O \rangle_1 + N \langle VO \rangle_1 + \frac{N^2}{2!} \langle V^2 O \rangle_1 + \frac{N^3}{3!} \langle V^3 O \rangle_1 + \dots}{1 + N \langle V \rangle_1 + \frac{N^2}{2!} \langle V^2 \rangle_1 + \frac{N^3}{3!} \langle V^3 \rangle_1 + \dots} \\ &= \langle O \rangle_1 \frac{1 + N \langle VO \rangle_1 / \langle O \rangle_1 + \frac{N^2}{2!} \langle V^2 O \rangle_1 / \langle O \rangle_1 + \frac{N^3}{3!} \langle V^3 O \rangle_1 / \langle O \rangle_1 + \dots}{1 + N \langle V \rangle_1 + \frac{N^2}{2!} \langle V^2 \rangle_1 + \frac{N^3}{3!} \langle V^3 \rangle_1 + \dots} \end{aligned} \quad (26)$$

where $\langle \cdot \rangle_1$ represents an average over the Gaussian fluctuation $\mathcal{N}(\mathbf{q}_1, -\frac{1}{N} \nabla^2 S(\mathbf{q}_1))$. We see that the series in the denominator contains terms of the form $\frac{N^k}{k!} \langle V^k \rangle_1$ while the numerator depends on terms of the form $\frac{N^k}{k!} \langle V^k O \rangle_1 / \langle O \rangle_1$. In either of these power series, the k -th term can contribute at most

$$\frac{N^k \langle V^k O \rangle_1}{\langle O \rangle_1}, N^k \langle V^k \rangle_1 \sim \begin{cases} \mathcal{O}(N^{-(k+1)-2}) & k \text{ odd} \\ \mathcal{O}(N^{-k-2}) & k \text{ even} \end{cases} \quad (27)$$

since V contributes only cubic and higher terms. Thus each term in the numerator and denominator's series contains increasing powers of $1/N$. Concretely, each of the two series have terms of order $\{N^0, N^{-1}, N^{-1}, N^{-2}, N^{-2}, \dots\}$. Thus any quantity of the form $\frac{\langle hO \rangle_1}{\langle hO \rangle_1}$ admits a ratio of power series in powers of $1/N$. One could truncate each of the series in the numerator and denominator to a desired order in N . Alternatively, the denominator could be expanded giving a single series (the cumulant expansion [56]). The first few terms in the cumulant expansion have the form

$$\begin{aligned} \langle O \rangle_N &= \langle O \rangle_1 + N[\langle OV \rangle_1 - \langle O \rangle_1 \langle V \rangle_1] \\ &\quad + \frac{N^2}{2} \langle V^2 O \rangle_1 - 2 \langle VO \rangle_1 \langle V \rangle_1 + 2 \langle V \rangle_1^2 \langle O \rangle_1 - \langle V^2 \rangle_1 \langle O \rangle_1 + \dots \end{aligned} \quad (28)$$

In this work, we mainly are interested in the leading order correction to $\langle O \rangle$ which can always be obtained with the truncation after the terms linear in V for any observable O .

D.1 Square Deviation from DMFT

We will now analyze the fluctuation statistics of our order parameters around the saddle point $(\mathbf{q} - \mathbf{q}_1)(\mathbf{q} - \mathbf{q}_1)^{\gt}_N$ which has the form

$$\begin{aligned} (\mathbf{q} - \mathbf{q}_1)(\mathbf{q} - \mathbf{q}_1)^{\gt}_N &= \frac{(\mathbf{q} - \mathbf{q}_1)(\mathbf{q} - \mathbf{q}_1)^{\gt}_1 + N \langle V \rangle (\mathbf{q} - \mathbf{q}_1)(\mathbf{q} - \mathbf{q}_1)^{\gt}_1 + \dots}{1 + N \langle V \rangle_1 + \dots} \\ &= \frac{\frac{1}{N} + \mathcal{O}(N^{-2})}{1 + \mathcal{O}(N^{-1})} \sim \frac{1}{N} + \mathcal{O}(N^{-2}), \end{aligned} \quad (29)$$

as stated in the main text and verified empirically in Figure 3 (a). The reason that the terms in the numerator involving V can be no larger than $\mathcal{O}(N^{-2})$ comes from vanishing of odd moments for $\mathbf{q} - \mathbf{q}_1$ in the unperturbed distribution. Thus the leading expression for $(\mathbf{q} - \mathbf{q}_1)(\mathbf{q} - \mathbf{q}_1)^{\gt}$ only depends on Σ and not on V .

D.2 Mean Deviation from DMFT

Although the square displacement from DMFT only depended on Σ and not on V , we note that the *average order parameter displacement* $\langle \mathbf{q} - \mathbf{q}_1 \rangle$ does receive a $\mathcal{O}(1/N)$ correction that depends on the perturbed potential V

$$\begin{aligned} \langle \mathbf{q} - \mathbf{q}_1 \rangle_N &= \frac{\langle \mathbf{q} - \mathbf{q}_1 \rangle_1 + N \langle (\mathbf{q} - \mathbf{q}_1) V \rangle_1 + \frac{N^2}{2} \langle (\mathbf{q} - \mathbf{q}_1) V^2 \rangle_1 + \dots}{1 + N \langle V \rangle_1 + \frac{N^2}{2} \langle V^2 \rangle_1 + \dots} \\ &\stackrel{D}{\sim} \frac{\frac{\partial V}{\partial \mathbf{q}}_1 + \mathcal{O}(N^{-2})}{1 + \mathcal{O}(N^{-1})} \stackrel{E}{\sim} \frac{\partial V}{\partial \mathbf{q}}_1 + \mathcal{O}(N^{-2}). \end{aligned} \quad (30)$$

where in the last line we used Stein's lemma (Gaussian integration by parts) for the Gaussian distribution over \mathbf{q} . Note that $\frac{\partial V}{\partial \mathbf{q}}_1 \sim \mathcal{O}\left(\frac{1}{N}\right)$ since the derivative of the cubic term in V gives a quadratic function of $\mathbf{q} - \mathbf{q}_1$, whose average must be $\mathcal{O}(N^{-1})$. In this work, we focus primarily on the structure of the propagator, but outline a general recipe for getting the leading mean correction in Appendix G and H.2.

D.3 Covariance of Order Parameters

Lastly, we combine the previous two observations to reason about the scaling of the order parameter covariance over initializations. We note that the leading covariance of the order parameters over random initializations is also given by the propagator: $\text{Cov}(\mathbf{q}) \sim \frac{1}{N} + \mathcal{O}(N^{-2})$, since

$$\begin{aligned} \text{Cov}(\mathbf{q}) &= \frac{D}{D} \frac{E}{E} (\mathbf{q} - \langle \mathbf{q} \rangle_N)(\mathbf{q} - \langle \mathbf{q} \rangle_N)^{\gt} \\ &= (\mathbf{q} - \mathbf{q}_1)(\mathbf{q} - \mathbf{q}_1)^{\gt}_N - (\mathbf{q}_1 - \langle \mathbf{q} \rangle_N)(\mathbf{q}_1 - \langle \mathbf{q} \rangle_N)^{\gt}_N \\ &\sim \frac{1}{N} + \mathcal{O}(N^{-2}) \end{aligned} \quad (31)$$

due to the arguments above which showed that $(\mathbf{q} - \mathbf{q}_1)(\mathbf{q} - \mathbf{q}_1)^{\gt} \sim \frac{1}{N} + \mathcal{O}(N^{-2})$ and that $\mathbf{q}_1 - \langle \mathbf{q} \rangle_N \sim \mathcal{O}(N^{-1})$. Therefore, in the leading order picture, it is safe to associate $\text{Cov}(\mathbf{q})$ with the covariance of order parameters over random initializations of the network weights.

E Propagator Structure for the full DMFT Action

In this section, we examine the propagator structure for the full DMFT action. This action is modified from other prior works [9, 46] to include the evolution of the network prediction errors $\Delta(t)$. Those prior works noted that Δ and the NTK K are deterministic functions of deterministic order parameters $\{\Phi, G\}$ in the $N \rightarrow \infty$ limit so those authors did not explicitly include Δ or K in the action. At finite width N , including Δ, K in the action is crucial as the fluctuation in prediction errors Δ has significant consequences for dynamical fluctuations of kernels through the preactivation and

pre-gradient fields. In this section, we will mainly focus on gradient flow, but we describe large step size in Appendix M.

$$\begin{aligned}
S = & \int_0^t dt ds \int \mathcal{D}h \hat{\Phi}(t, s) \Phi(t, s) + \hat{G}(t, s) G(t, s) - \gamma^2 A(s, t) B(t, s) \\
& + \int_0^t dt \hat{\Delta}(t) \Delta(t) - y + \int_0^t ds \Theta(t-s) K(s) \Delta(s) \\
& + \int_0^t dt \hat{K}(t) K(t) - \int_0^t G^{+1}(t) \Phi(t) \\
& + \ln \mathcal{Z}[\hat{\Phi}, \hat{G}, \hat{\Delta}, \hat{K}, G^{+1}, A^{-1}, B]
\end{aligned} \tag{32}$$

where the single site moment generating functionals \mathcal{Z} have the form

$$\begin{aligned}
\mathcal{Z} &= \mathbb{E}_{\mathcal{P}(h(t), z(t), g)} \exp \left[- \int_0^t dt ds \int \mathcal{D}h \phi(h(t)) \phi(h(s)) \hat{\Phi}(t, s) + g(t) g(s) \hat{G}(t, s) \right] \\
h(t) &= u(t) + \gamma \int_0^t ds \int \mathcal{D}h \Phi^{-1}(t, s) \Delta(s) + A^{-1}(t, s) g(s), \{u(t)\} \sim \mathcal{GP}(0, \Sigma^{-1}) \\
z(t) &= r(t) + \gamma \int_0^t ds \int \mathcal{D}h G^{+1}(t, s) \Delta(s) + B(t, s) \phi(h(s)), \{r(t)\} \sim \mathcal{GP}(0, \mathbf{G}^{+1})
\end{aligned} \tag{33}$$

with $g(t) = \dot{\phi}(h(t)) z(t)$. The saddle point equations give the infinite width evolution of our order parameters.

$$\begin{aligned}
\frac{\partial S}{\partial \hat{\Phi}(t, s)} &= \Phi(t, s) - \phi(h(t)) \phi(h(s)) = 0 \\
\frac{\partial S}{\partial \hat{G}(t, s)} &= G(t, s) - g(t) g(s) = 0 \\
\frac{\partial S}{\partial A(s, t)} &= -\gamma^2 B(t, s) + \gamma \frac{\partial \phi(h(t))}{\partial r(s)} = 0 \\
\frac{\partial S}{\partial B(s, t)} &= -\gamma^2 A(t, s) + \gamma \frac{\partial g(t)}{\partial u(s)} = 0 \\
\frac{\partial S}{\partial \hat{K}(t)} &= K(t) - \int_0^t G^{+1}(t, t) \Phi(t, t) = 0 \\
\frac{\partial S}{\partial \hat{\Delta}(t)} &= \Delta(t) - y + \int_0^t ds \int \mathcal{D}h K(s) \Delta(s) = 0
\end{aligned} \tag{34}$$

These equations exactly recover the mean field description obtained [9]. Note that $\langle \cdot \rangle$ for field averages is an average defined by \mathcal{Z} and is distinct from the types averages $\langle \cdot \rangle, \langle \cdot \rangle_{\mathcal{P}}$ we have been considering over the order parameters \mathbf{q} . The complementary set of equations for the primal variables, such as $\frac{\partial S}{\partial \hat{\Phi}(t, s)} = 0$, give that $\hat{K} = \hat{\Delta} = \hat{\Phi} = \hat{G} = 0$ at the saddle point. We now set out to compute the Hessian $\nabla_{\mathbf{q}}^2 S$. To simplify the set of expressions, we will only explicitly write out the nonvanishing

blocks. We will start with second derivatives involving only pairs of dual variables $\{\hat{\Phi}, \hat{G}, A, B\}$

$$\begin{aligned}
\frac{\partial^2 S}{\partial \hat{\Phi}(t, s) \partial \hat{\Phi}(t^\theta, s^\theta)} &= \phi(h(t))\phi(h(s))\phi(h(t^\theta))\phi(h(s^\theta)) - \Phi(t, s)\Phi(t^\theta, s^\theta) \\
&\equiv \kappa(t, s, t^\theta, s^\theta) \\
\frac{\partial^2 S}{\partial \hat{G}(t, s) \partial \hat{G}(t^\theta, s^\theta)} &= g(t)g(s)g(t^\theta)g(s^\theta) - G(t, s)G(t^\theta, s^\theta) \\
&\equiv \kappa^G(t, s, t^\theta, s^\theta) \\
\frac{\partial^2 S}{\partial \hat{\Phi}(t, s) \partial \hat{G}(t^\theta, s^\theta)} &= \phi(h(t))\phi(h(s))g(t^\theta)g(s^\theta) - \Phi(t, s)G(t^\theta, s^\theta) \\
&\equiv \kappa_*^G(t, s, t^\theta, s^\theta) \\
\frac{\partial^2 S}{\partial \hat{\Phi}(t, s) \partial A^{-1}(s^\theta, t^\theta)} &= -\gamma \frac{\partial \phi(h(t))}{\partial u(s^\theta)} \phi(h(s))g(t^\theta) \\
&\quad - \gamma \phi(h(t)) \frac{\partial \phi(h(t))}{\partial u(s^\theta)} g(t^\theta) \\
&\quad - \gamma \phi(h(t))\phi(h(s)) \frac{\partial g(t^\theta)}{\partial u(s^\theta)} - \gamma^2 \Phi(t, s)B^{-1}(t^\theta, s^\theta) \\
&\equiv -\gamma \kappa_*^{B^{-1}}(t, s) \\
\frac{\partial^2 S}{\partial \hat{\Phi}(t, s) \partial B^{-1}(s^\theta, t^\theta)} &= -\gamma \frac{\partial \phi(h(t))}{\partial r(s^\theta)} \phi(h(s))\phi(h(t^\theta)) \\
&\quad - \gamma \phi(h(t)) \frac{\partial \phi(h(t))}{\partial r(s^\theta)} \phi(h(t^\theta)) \\
&\quad - \gamma \phi(h(t))\phi(h(s)) \frac{\partial \phi(h(t^\theta))}{\partial r(s^\theta)} - \gamma^2 \Phi(t, s)A^{-1}(t^\theta, s^\theta) \\
&\equiv -\gamma \kappa_*^A(t, s) \\
\frac{\partial^2 S}{\partial \hat{G}(t, s) \partial A^{-1}(s^\theta, t^\theta)} &= -\gamma \frac{\partial g(t)}{\partial u(s^\theta)} g(s)g(t^\theta) - \gamma g(t) \frac{\partial g(t)}{\partial u(s^\theta)} g(t^\theta) \\
&\quad - \gamma g(t)g(s) \frac{\partial g(t^\theta)}{\partial u(s^\theta)} - \gamma^2 G(t, s)B^{-1}(t^\theta, s^\theta) \\
&\equiv -\gamma \kappa_*^{G B^{-1}}(t, s) \\
\frac{\partial^2 S}{\partial \hat{G}(t, s) \partial B^{-1}(s^\theta, t^\theta)} &= -\gamma \frac{\partial g(t)}{\partial r(s^\theta)} g(s)\phi(h(t^\theta)) - \gamma g(t) \frac{\partial g(t)}{\partial r(s^\theta)} \phi(h(t^\theta)) \\
&\quad - \gamma g(t)g(s) \frac{\partial \phi(h(t^\theta))}{\partial r(s^\theta)} - \gamma^2 G(t, s)A^{-1}(t^\theta, s^\theta) \\
&\equiv -\gamma \kappa_*^{G A^{-1}}(t, s) \\
\frac{\partial^2 S}{\partial A^{-1}(t, s) \partial B^{-1}(s^\theta, t^\theta)} &= -\gamma^2 \delta \delta(t - t^\theta) \delta(s - s^\theta) \\
\frac{\partial^2 S}{\partial A^{-1}(s, t) \partial B^{-1}(s^\theta, t^\theta)} &= \gamma^2 \frac{\partial^2}{\partial u(s) \partial r(s^\theta)} g(t)\phi(h(t^\theta)) - \gamma^4 B^{-1}(t, s)A^{-1}(t^\theta, s^\theta) \\
&\equiv \kappa^{B^{-1} A^{-1}}(t, s, t^\theta, s^\theta)
\end{aligned} \tag{35}$$

Next, we consider the second derivatives involving only primal variables $\{\Phi^{\cdot}, G^{\cdot}, K^{\cdot}, \Delta\}$ which all vanish

$$\begin{aligned}
\frac{\partial^2 S}{\partial \Phi^{\cdot}(t, s) \partial \Phi^{\cdot\theta}(t^{\theta}, s^{\theta})} &= 0 \\
\frac{\partial^2 S}{\partial G^{\cdot}(t, s) \partial G^{\cdot\theta}(t^{\theta}, s^{\theta})} &= 0 \\
\frac{\partial^2 S}{\partial \Phi^{\cdot}(t, s) \partial G^{\cdot\theta}(t^{\theta}, s^{\theta})} &= 0 \\
\frac{\partial^2 S}{\partial \Phi^{\cdot}(t, s) \partial K^{\cdot}(s^{\theta})} &= 0 \\
\frac{\partial^2 S}{\partial G^{\cdot}(t, s) \partial K^{\cdot}(s^{\theta})} &= 0 \\
\frac{\partial^2 S}{\partial \Phi^{\cdot}(t, s) \partial \Delta^{\cdot}(s^{\theta})} &= 0 \\
\frac{\partial^2 S}{\partial G^{\cdot}(t, s) \partial \Delta^{\cdot}(s^{\theta})} &= 0 \\
\frac{\partial^2 S}{\partial K^{\cdot}(t) \partial K^{\cdot}(s)} &= 0 \\
\frac{\partial^2 S}{\partial K^{\cdot}(t) \partial \Delta^{\cdot}(s)} &= 0 \\
\frac{\partial^2 S}{\partial \Delta^{\cdot}(t) \partial \Delta^{\cdot}(s)} &= 0
\end{aligned} \tag{36}$$

Now we consider all derivatives which involve one of the dual variables $\{\hat{\Phi}^{\cdot}, \hat{G}^{\cdot}, A^{\cdot}, B^{\cdot}\}$ and the primal variable Δ

$$\begin{aligned}
\frac{\partial^2 S}{\partial \hat{\Phi}^{\cdot}(t, s) \partial \Delta^{\cdot}(t^{\theta})} &= - \frac{\partial}{\partial \Delta^{\cdot}(t^{\theta})} [\phi(h^{\cdot}(t)) \phi(h^{\cdot}(s))] \equiv -D^{\cdot}(t, s, t^{\theta}) \\
\frac{\partial^2 S}{\partial \hat{G}^{\cdot}(t, s) \partial \Delta^{\cdot}(t^{\theta})} &= - \frac{\partial}{\partial \Delta^{\cdot}(t^{\theta})} [g^{\cdot}(t) g^{\cdot}(s)] \equiv -D^{G^{\cdot}}(t, s, t^{\theta}) \\
\frac{\partial^2 S}{\partial A^{\cdot-1}(s, t) \partial \Delta^{\cdot}(t^{\theta})} &= \gamma \frac{\partial}{\partial \Delta^{\cdot}(t^{\theta}) \partial u^{\cdot}(s)} g^{\cdot}(t) \equiv \gamma D^{B^{\cdot-1}}(t, s, t^{\theta}) \\
\frac{\partial^2 S}{\partial B^{\cdot}(s, t) \partial \Delta^{\cdot}(t^{\theta})} &= \gamma \frac{\partial}{\partial \Delta^{\cdot}(t^{\theta}) \partial r^{\cdot}(s)} \phi(h^{\cdot}(t)) \equiv \gamma D^{A^{\cdot}}(t, s, t^{\theta})
\end{aligned}$$

Now, we consider the second derivatives involving one derivative on a dual variable $\{\hat{\Phi}, \hat{G}, A, B\}$ and one of the primal variables $\{\Phi, G\}$.

$$\begin{aligned}
\frac{\partial^2 S}{\partial \hat{\Phi}(t, s) \partial \Phi(t^\theta, s^\theta)} &= \delta_{:,0} \delta_{:,0} \delta(t - t^\theta) \delta(s - s^\theta) \\
&\quad - \delta_{:,1;0} \frac{\partial}{\partial \Phi^{-1}(t^\theta, s^\theta)} \phi(h(t)) \phi(h(s)) \\
&\equiv \delta_{:,0} \delta_{:,0} \delta(t - t^\theta) \delta(s - s^\theta) - \delta_{:,1;0} D^{\Phi^{-1}}(t, s, t^\theta, s^\theta) \\
\frac{\partial^2 S}{\partial \hat{G}(t, s) \partial G(t^\theta, s^\theta)} &= \delta_{:,0} \delta_{:,0} \delta(t - t^\theta) \delta(s - s^\theta) - \delta_{:,+1;0} \frac{\partial}{\partial G^{+1}(t^\theta, s^\theta)} g(t) g(s) \\
&\equiv \delta_{:,0} \delta_{:,0} \delta(t - t^\theta) \delta(s - s^\theta) - \delta_{:,+1;0} D^{G^{+1}}(t, s, t^\theta, s^\theta) \\
\frac{\partial^2 S}{\partial \hat{\Phi}(t, s) \partial G^{+1}(t^\theta, s^\theta)} &= - \frac{\partial}{\partial G^{+1}(t^\theta, s^\theta)} \phi(h(t)) \phi(h(s)) \equiv -D^{\Phi^{-1};G^{+1}}(t, s, t^\theta, s^\theta) \\
\frac{\partial^2 S}{\partial \hat{G}(t, s) \partial \Phi^{-1}(t^\theta, s^\theta)} &= - \frac{\partial}{\partial \Phi^{-1}(t^\theta, s^\theta)} g(t) g(s) \equiv -D^{G^{+1};\Phi^{-1}}(t, s, t^\theta, s^\theta) \\
\frac{\partial^2 S}{\partial A^{-1}(s, t) \partial \Phi^{-1}(t^\theta, s^\theta)} &= \gamma \frac{\partial}{\partial \Phi^{-1}(t^\theta, s^\theta)} \frac{\partial g(t)}{\partial r(s)} \equiv \gamma D^{B^{-1};\Phi^{-1}}(t, s, t^\theta, s^\theta) \\
\frac{\partial^2 S}{\partial B(s, t) \partial \Phi^{-1}(t^\theta, s^\theta)} &= \gamma \frac{\partial}{\partial \Phi^{-1}(t^\theta, s^\theta)} \frac{\partial \phi(h(t))}{\partial u(s)} \equiv \gamma D^{A^{-1};\Phi^{-1}}(t, s, t^\theta, s^\theta) \\
\frac{\partial^2 S}{\partial A^{-1}(s, t) \partial G^{+1}(t^\theta, s^\theta)} &= \gamma \frac{\partial}{\partial G^{+1}(t^\theta, s^\theta)} \frac{\partial g(t)}{\partial r(s)} \equiv \gamma D^{B^{-1};G^{+1}}(t, s, t^\theta, s^\theta) \\
\frac{\partial^2 S}{\partial B(s, t) \partial G^{+1}(t^\theta, s^\theta)} &= \gamma \frac{\partial}{\partial G^{+1}(t^\theta, s^\theta)} \frac{\partial \phi(h(t))}{\partial u(s)} \equiv \gamma D^{A^{-1};G^{+1}}(t, s, t^\theta, s^\theta) \quad (37)
\end{aligned}$$

We note that terms such as $\frac{\partial}{\partial \Phi^{-1}(t^\theta, s^\theta)} \phi(h(t)) \phi(h(s))$ can be further decomposed since the average over the $\{u(t)\} \sim \mathcal{GP}(0, \Sigma^{-1})$ and h 's explicit dynamics both depend on Φ^{-1}

$$\begin{aligned}
\frac{\partial}{\partial \Phi^{-1}(t^\theta, s^\theta)} \phi(h(t)) \phi(h(s)) &= \frac{1}{2} \frac{\partial^2}{\partial u(t^\theta) \partial u(s^\theta)} \phi(h(t)) \phi(h(s)) \\
&\quad + \frac{\partial}{\partial \Phi^{-1}(t^\theta, s^\theta)} \phi(h(t)) \phi(h(s)) \quad (38)
\end{aligned}$$

where the first term comes from differentiating the Gaussian probability density for u (e.g. Price's theorem) and the second term is an explicit derivative of the preactivation fields with u treated as constant. Next we consider the nonvanishing terms which involve $\{\hat{\Delta}, \hat{K}, \Delta, K\}$ which give

$$\begin{aligned}
\frac{\partial^2 S}{\partial \hat{\Delta}(t) \partial \Delta(s)} &= \delta_{:,0} \delta(t - s) + \Theta(t - s) K(s) \\
\frac{\partial^2 S}{\partial \hat{\Delta}(t) \partial K(s)} &= \delta_{:,0} \Theta(t - s) \Delta(s) \\
\frac{\partial^2 S}{\partial \hat{K}(t) \partial K(t^\theta)} &= \delta_{:,0} \delta_{:,0} \delta(t - t^\theta) \\
\frac{\partial^2 S}{\partial \hat{K}(t) \partial \Phi^{-1}(t^\theta, s^\theta)} &= \delta_{:,0} \delta_{:,0} G^{+1}(t^\theta, s^\theta) \delta(t - t^\theta) \delta(t - s^\theta) \\
\frac{\partial^2 S}{\partial \hat{K}(t) \partial G^{+1}(t^\theta, s^\theta)} &= \delta_{:,0} \delta_{:,0} \Phi^{-1}(t^\theta, s^\theta) \delta(t - t^\theta) \delta(t - s^\theta) \quad (39)
\end{aligned}$$

This enumerates all possible non-vanishing terms in the Hessian. We can now construct a block matrix of these Hessians by partitioning our order parameters $\mathbf{q} = [\mathbf{q}_1, \mathbf{q}_2]^>$ where

$$\mathbf{q}_1 = \text{Vec}\{\Phi^{\cdot}(t, s), G^{\cdot}(t, s), K^{\cdot}(t), \Delta^{\cdot}(t), \hat{\Phi}^{\cdot}(t, s), \hat{G}^{\cdot}(t, s), \hat{K}^{\cdot}(t), \hat{\Delta}^{\cdot}(t)\} \quad (40)$$

$$\mathbf{q}_2 = \text{Vec}\{A^{\cdot}(t, s), B^{\cdot}(t, s)\}. \quad (41)$$

This choice will become apparent shortly.

$$\nabla_{\mathbf{q}}^2 S = \begin{pmatrix} \nabla_{\mathbf{q}_1}^2 S & \nabla_{\mathbf{q}_1 \mathbf{q}_2}^2 S \\ \nabla_{\mathbf{q}_2 \mathbf{q}_1}^2 S & \nabla_{\mathbf{q}_2}^2 S \end{pmatrix} \quad (42)$$

To calculate the full propagator $\mathbb{G} = -\nabla_{\mathbf{q}}^2 S^{-1}$, we will assume invertibility of the upper block $\mathbb{G}^0 = -\nabla_{\mathbf{q}_1}^2 S^{-1}$ and use this in the Schur complement

$$\begin{aligned} \mathbb{G} &= -\nabla_{\mathbf{q}}^2 S^{-1} = \begin{pmatrix} 11 & 12 \\ 21 & 22 \end{pmatrix} \\ 11 &= \mathbb{G}^0 - \mathbb{G}^0 \nabla_{\mathbf{q}_1 \mathbf{q}_2}^2 S \nabla_{\mathbf{q}_2}^2 S + (\nabla_{\mathbf{q}_2 \mathbf{q}_1}^2 S)^0 (\nabla_{\mathbf{q}_1 \mathbf{q}_2}^2 S)^{-1} \nabla_{\mathbf{q}_2 \mathbf{q}_1}^2 S \mathbb{G}^0 \\ 12 &= \mathbb{G}^0 \nabla_{\mathbf{q}_1 \mathbf{q}_2}^2 S \nabla_{\mathbf{q}_2}^2 S + (\nabla_{\mathbf{q}_2 \mathbf{q}_1}^2 S)^0 (\nabla_{\mathbf{q}_1 \mathbf{q}_2}^2 S)^{-1} \\ 22 &= -\nabla_{\mathbf{q}_2}^2 S + (\nabla_{\mathbf{q}_2 \mathbf{q}_1}^2 S)^0 (\nabla_{\mathbf{q}_1 \mathbf{q}_2}^2 S)^{-1} \end{aligned} \quad (43)$$

We now need to solve for $\mathbb{G}^0 = -\nabla_{\mathbf{q}_1}^2 S^{-1}$. To perform this inverse, we again partition \mathbf{q}_1 into two sets of order parameters $\mathbf{q}_1 = [\mathbf{q}_1^1, \mathbf{q}_1^2]$ where $\mathbf{q}_1^1 = \text{Vec}\{\Phi^{\cdot}(t, s), G^{\cdot}(t, s), K^{\cdot}(t), \Delta^{\cdot}(t)\}$ and $\mathbf{q}_1^2 = \text{Vec}\{\hat{\Phi}^{\cdot}(t, s), \hat{G}^{\cdot}(t, s), \hat{K}^{\cdot}(t), \hat{\Delta}^{\cdot}(t)\}$

$$\nabla_{\mathbf{q}_1}^2 S = \begin{pmatrix} \mathbf{0} & \mathbf{U}^> \\ \mathbf{U} & \end{pmatrix}, \quad \mathbb{G}^0 \equiv \nabla_{\mathbf{q}_1^1}^2 S, \quad \mathbf{U} \equiv \nabla_{\mathbf{q}_1^2 \mathbf{q}_1^1}^2 S \quad (44)$$

We seek a physically sensible inverse where the variance of \mathbf{q}_1^2 is vanishing [51, 53]. This leads to the following sub-propagator \mathbb{G}^0

$$\mathbb{G}^0 = -[\nabla_{\mathbf{q}_1^1}^2 S]^{-1} = \begin{pmatrix} \mathbf{U}^{-1} & [\mathbf{U}^{-1}]^> \\ -[\mathbf{U}^>]^{-1} & \mathbf{0} \end{pmatrix} \quad (45)$$

Thus given \mathbf{U} , we can solve for \mathbb{G}^0 and ultimately for the full propagator \mathbb{G} . The relevant entries in \mathbb{G}^0 and \mathbf{U} are given by those second derivatives calculated above. We note that each of the field derivatives needed for \mathbf{U} can be computed implicitly from the field dynamics. For example, for the $\Delta^{\cdot}(t)$ derivatives we have

$$\begin{aligned} \frac{\partial}{\partial \Delta^{\cdot}(t^\theta)} h^{\cdot}(t) &= \gamma \Theta(t - t^\theta) \Phi^{\cdot-1}(t, t^\theta) g^{\cdot}(t^\theta) \\ &\quad + \gamma \int_0^t ds \times A^{\cdot-1}(t, s) + \Phi^{\cdot-1}(t, s) \Delta^{\cdot}(s) \frac{\partial g^{\cdot}(s)}{\partial \Delta^{\cdot}(t^\theta)} \\ \frac{\partial}{\partial \Delta^{\cdot}(t^\theta)} z^{\cdot}(t) &= \gamma \Theta(t - t^\theta) G^{\cdot+1}(t, t^\theta) \phi(h^{\cdot}(t^\theta)) \\ &\quad + \gamma \int_0^t ds \times B^{\cdot}(t, s) + G^{\cdot+1}(t, s) \Delta^{\cdot}(s) \frac{\partial \phi(h^{\cdot}(s))}{\partial \Delta^{\cdot}(t^\theta)} \end{aligned} \quad (46)$$

These can then be used in the averages such as $\frac{\partial}{\partial \Delta^{\cdot}(t^\theta)} \phi(h^{\cdot}(t)) \phi(h^{\cdot}(s))$. Similarly, we can compute terms such as $\frac{\partial h^{\cdot}(t)}{\partial \Phi^{\cdot-1}(t^\theta, s^\theta)}$ through the following closed equations

$$\begin{aligned} \frac{\partial h^{\cdot}(t)}{\partial \Phi^{\cdot-1}(t^\theta, s^\theta)} &= \gamma \delta(t - t^\theta) \delta^{\cdot} \Theta(t - s^\theta) \Delta^{\cdot}(s^\theta) \\ &\quad + \gamma \int_0^t ds \times A^{\cdot-1}(t, s) + \Delta^{\cdot}(s) \Phi^{\cdot-1}(t, s) \frac{\partial g^{\cdot}(s)}{\partial \Phi^{\cdot-1}(t^\theta, s^\theta)} \\ \frac{\partial z^{\cdot}(t)}{\partial \Phi^{\cdot-1}(t^\theta, s^\theta)} &= \gamma \int_0^t ds \times B^{\cdot}(t, s) + \Delta^{\cdot}(s) G^{\cdot+1}(t, s) \frac{\partial \phi(h^{\cdot}(s))}{\partial \Phi^{\cdot-1}(t^\theta, s^\theta)} \end{aligned} \quad (47)$$

These terms can then be used to compute quantities like D^{\cdot} .

F Solving for the Propagator

In this section we sketch out the required steps to obtain the propagator Σ .

- Step 1: Solve the infinite width DMFT equations for q_1 which include the prediction error dynamics $\Delta(t)$, the feature kernels $\Phi(t, s)$, gradient kernels $G(t, s)$. This step corresponds to algorithm in Bordelon & Pehlevan '22 and defines the dynamics one would expect at infinite width [9]. See below for more detail.
- Step 2: Compute the entries of the Hessian of S evaluated at the q_1 computed in the first step. Some of these entries look like fourth cumulants of features like $\kappa = \langle \phi(h)^4 \rangle - 3 \langle \phi(h)^2 \rangle^2$ and some of them measure sensitivity of one order parameter to a perturbation in another order parameter $D = \frac{\partial \langle \phi(h) \rangle}{\partial q_1} \langle \phi(h) \rangle^2$. The averages $\langle \cdot \rangle$ used to calculate κ and D should be performed over the infinite width stochastic processes for preactivations h which are defined in equation (19).
- Step 3: After populating the entries of the block matrix for the Hessian $\nabla^2 S$, we then calculate the propagator Σ with a matrix inversion. Since we discretized time, this is a finite dimensional matrix.

The step 1 above demands a solution to the infinite width DMFT equations (solving for the saddle point q_1). We will now give a detailed set of instructions about how the infinite width limit for q_1 is solved (step 1 above). This corresponds to the algorithm of Bordelon & Pehlevan 2022 to solve the saddle point equations $\frac{\partial}{\partial q} S(q)|_{q_1} = 0$ [9].

- Step 1: Start with a guess for the kernels $\Phi(t, s), G(t, s)$ and for the predictions through time $f(t)$. We usually use the lazy limit (e.g. $\Phi(t, s) = \Phi(0, 0) \dots$) as an initial guess.
- Step 2: Sample Gaussian sources $u(t)$ and $r(t)$ based on the current covariances Φ and G .
- Step 3: For each sample, solve integral equations for $h(t)$ and $z(t)$.

$$\begin{aligned} h(t) &= u(t) + \gamma \int_0^t ds [A^{-1}(t, s) + \Phi^{-1}(t, s)] [\phi(h(s))z(s)] \\ z(t) &= r(t) + \gamma \int_0^t ds [B^{-1}(t, s) + G^{-1}(t, s)] \phi(h(s)) \end{aligned} \quad (48)$$

These will be samples from the single site distribution for h, z

- Step 4: Average over the Monte Carlo samples to produce a new estimate of the kernels: $\Phi(t, s) = \langle \phi(h(t))\phi(h(s)) \rangle$. A similar procedure is performed for G and the response functions A, B .
- Step 5: Compute the NTK estimate $K(t) = \langle G^{-1}(t, t)\Phi(t, t) \rangle$ and then integrate prediction dynamics from the dynamics of the NTK $\frac{d}{dt}f(t) = K(t)\Delta(t)$.
- Repeat steps 2-5 until the order parameters converge.

Below we provide a pseudocode algorithm to solve for the propagator elements.

Algorithm 1: Propagator Solver

Data: K^x, \mathbf{y} , Initial Guesses $\{\Phi, G\}_{L=1}^L, \{A, B\}_{L=1}^L$, Sample count \mathcal{S} , Update Speed β

Result: Propagator Matrix

- 1 Solve DMFT equations with Algorithm 2 for order parameters $f(t), \Phi(t, s), \dots$;
 - 2 Draw \mathcal{S} samples $\{u_{:,n}(t)\}_{n=1}^{\mathcal{S}} \sim \mathcal{GP}(0, \Phi^{-1}), \{r_{:,n}(t)\}_{n=1}^{\mathcal{S}} \sim \mathcal{GP}(0, G^{-1})$;
 - 3 Integrate dynamics for each sample to get $\{h_{:,n}(t), z_{:,n}(t)\}_{n=1}^{\mathcal{S}}$;
 - 4 Estimate κ functions with Monte Carlo integration, for instance
 - 5 $\kappa_{\text{P}}(t, s, t^\theta, s^\theta) = \frac{1}{\mathcal{S}} \sum_{n \in [S]} \langle \phi(h_{:,n}(t))\phi(h_{:,n}(s))\phi(h_{:,n}(t^\theta))\phi(h_{:,n}(s^\theta)) \rangle - \Phi(t, s)\Phi(t^\theta, s^\theta)$;
 - 6 For each sample, compute field sensitivities to error signals, such as $\frac{\partial h_{:,n}(t)}{\partial (t^\theta, s^\theta)}$, and kernels $\frac{\partial h_{:,n}(t)}{\partial (t^\theta, s^\theta)}$ implicitly using equations (46) (47) ;
 - 7 Use these sensitivities to compute the necessary D tensors such as $D = \frac{1}{\mathcal{S}} \sum_{n \in [S]} \frac{\partial \langle \phi(h_{:,n}(t))\phi(h_{:,n}(s)) \rangle}{\partial (t^\theta, s^\theta)}$;
 - 8 Invert U matrix and compute Σ in equation (45);
 - 9 Compute the Schur-complement in equation (43) to handle the response functions ;
-

The above propagator solver builds on the solution to the DMFT equations which is provided below.

Algorithm 2: Alternating Monte Carlo Solution to Saddle Point Equations

Data: K^x, \mathbf{y} , Initial Guesses $\{\tilde{\mathbf{K}}, \tilde{\mathbf{G}}\}_{l=1}^L, \{\tilde{\mathbf{A}}, \tilde{\mathbf{B}}\}_{l=1}^L, \beta$, Sample count \mathcal{S} , Update Speed β
Result: Final Kernels $\{\tilde{\mathbf{K}}, \tilde{\mathbf{G}}\}_{l=1}^L, \{\tilde{\mathbf{A}}, \tilde{\mathbf{B}}\}_{l=1}^L$, Network predictions through training $f^*(t)$

```

1  $\mathbf{K}^0 = K^x \otimes \mathbf{1}^{\mathcal{S}}, \mathbf{G}^{L+1} = \mathbf{1}^{\mathcal{S}};$ 
2 while Kernels Not Converged do
3   From  $\{\tilde{\mathbf{K}}, \tilde{\mathbf{G}}\}$  compute  $K^{NTK}(t, t)$  and solve  $\frac{d}{dt}f^*(t) = \mathbf{P} \Delta(t)K^{NTK}(t, t);$ 
4    $\ell = 1;$ 
5   while  $\ell < L + 1$  do
6     Draw  $\mathcal{S}$  samples  $\{u_{:,n}(t)\}_{n=1}^{\mathcal{S}} \sim \mathcal{GP}(0, \tilde{\mathbf{K}}^{-1}), \{r_{:,n}(t)\}_{n=1}^{\mathcal{S}} \sim \mathcal{GP}(0, \tilde{\mathbf{G}}^{-1});$ 
7     Integrate dynamics for each sample to get  $\{h_{:,n}(t), z_{:,n}(t)\}_{n=1}^{\mathcal{S}};$ 
8     Compute new  $\mathbf{P}, \tilde{\mathbf{G}}$  estimates:
9      $\tilde{\Phi}^*(t, s) = \frac{1}{\mathcal{S}} \sum_{n \in [S]} \phi(h_{:,n}(t))\phi(h_{:,n}(s)), \tilde{\mathbf{G}}^*(t, s) = \frac{1}{\mathcal{S}} \sum_{n \in [S]} g_{:,n}(t)g_{:,n}(s);$ 
10    Solve for Jacobians on each sample  $\frac{\partial \langle h_n \rangle}{\partial r_n^>}, \frac{\partial \langle g_n \rangle}{\partial u_n^>};$ 
11    Compute new  $\tilde{\mathbf{A}}, \tilde{\mathbf{B}}^{-1}$  estimates:
12     $\tilde{\mathbf{A}}^* = \frac{1}{\mathcal{S}} \sum_{n \in [S]} \frac{\partial \langle h_n \rangle}{\partial r_n^>}, \tilde{\mathbf{B}}^{-1} = \frac{1}{\mathcal{S}} \sum_{n \in [S]} \frac{\partial \langle g_n \rangle}{\partial u_n^>};$ 
13     $\ell \leftarrow \ell + 1;$ 
14  end
15   $\ell = 1;$ 
16  while  $\ell < L + 1$  do
17    Update feature kernels:  $\tilde{\mathbf{K}} \leftarrow (1 - \beta)\tilde{\mathbf{K}} + \beta\tilde{\mathbf{K}}^*, \tilde{\mathbf{G}} \leftarrow (1 - \beta)\tilde{\mathbf{G}} + \beta\tilde{\mathbf{G}}^*;$ 
18    if  $\ell < L$  then
19      Update  $\tilde{\mathbf{A}} \leftarrow (1 - \beta)\tilde{\mathbf{A}} + \beta\tilde{\mathbf{A}}^*, \tilde{\mathbf{B}} \leftarrow (1 - \beta)\tilde{\mathbf{B}} + \beta\tilde{\mathbf{B}}^*;$ 
20    end
21     $\ell \leftarrow \ell + 1$ 
22  end
23 end
24 return  $\{\tilde{\mathbf{K}}, \tilde{\mathbf{G}}\}_{l=1}^L, \{\tilde{\mathbf{A}}, \tilde{\mathbf{B}}\}_{l=1}^L, \{f^*(t)\}_{l=1}^L$ 

```

G Leading Correction to the Mean Order Parameters

In this section we use the propagator structure derived in the last section to reason about the leading finite size correction to $\langle \mathbf{q} \rangle$ at width N . Letting the indices i, j, k, n enumerate all entries of the order parameters in \mathbf{q} (technically this is a sum over samples and an integral over time for gradient flow), we find the leading Pade Approximant for the mean has the form (App D)

$$\begin{aligned} \langle q_i - q_i^1 \rangle_N &= \frac{N \langle (q_i - q_i^1) V \rangle_1 + \frac{N^2}{2} \langle (q_i - q_i^1) V^2 \rangle_1 \dots}{1 + N \langle V \rangle_1 + \frac{N^2}{2} \langle V^2 \rangle_1 + \dots} \\ &\sim \frac{1}{3!N} \times \frac{\partial^3 S}{\partial q_j \partial q_k \partial q_l} \langle \delta_i \delta_j \delta_k \delta_l \rangle_1 + \mathcal{O}(N^{-2}). \end{aligned} \quad (49)$$

$$= \frac{1}{2N} \times \frac{\partial^3 S}{\partial q_j \partial q_k \partial q_l} \Sigma_{ij} \Sigma_{kl} + \mathcal{O}(N^{-2}) \quad (50)$$

where $\delta_j = \sqrt{N}(q_j - q_j^1)$ and the derivatives are computed at the saddle point. In the last line, we utilized Wick's theorem and the permutation symmetry of the third derivative $\frac{\partial^3 S}{\partial q_i \partial q_j \partial q_k}$ to evaluate the four point averages in terms of the propagator Σ_{ij} , which was provided in the preceding section E. In practice computing even the full set of second derivatives for the DMFT action to get Σ is quite challenging. Despite the challenge of computing the mean order parameter correction, these corrections are relevant in practice and crucially distinguish the training timescales of deep networks at different widths as we show in Figures 7 and A.4.

G.1 Correction to Mean Predictions and Full MSE Correction

Supposing that we solved for the propagator $\Delta(t)$, using the formalism in the preceding section, we can compute the $\mathcal{O}(N^{-1})$ correction to the average network prediction error due to finite size. We let $\langle \Delta(t) \rangle$ represent the average of errors over an ensemble of width N networks.

$$\begin{aligned} \frac{d}{dt} \langle \Delta(t) \rangle &= - \langle K(t) \Delta(t) \rangle \\ &= - \langle K(t) \rangle \langle \Delta(t) \rangle - \text{Cov}(K(t), \Delta(t)) \\ &\sim - \langle K(t) \rangle \langle \Delta(t) \rangle - \frac{1}{N} \Sigma^K(t, t) + \mathcal{O}(N^{-2}) \end{aligned} \quad (51)$$

where $\Sigma^K(t, t)$ is the leading covariance (propagator element) between the kernel $K(t)$ and prediction error $\Delta(t)$. We see that the average kernel $\langle K(t) \rangle$ (which depends on the finite width N) plays an important role in characterizing the timescales of the average prediction dynamics. Once this equation is solved for $\langle \Delta(t) \rangle$, the square loss at width N and time t has the form

$$\langle \Delta(t) \rangle^2 \sim \left(1 - \frac{2}{N} \langle K(t) \rangle \right) \langle \Delta(t) \rangle^2 + \frac{2}{N} \langle \Delta(t) \rangle \langle \Delta^T(t) \rangle + \frac{1}{N} \Sigma(t, t) + \mathcal{O}(N^{-2}) \quad (52)$$

We will now comment on the structure of the cross term in this above solution. First, if $\langle K \rangle \succeq K^T$ and Σ^K is negligible then the average errors at finite width will decay more rapidly than the infinite width model. However, we suspect that in general, $\langle K \rangle - K^T$ contains many negative eigenvalues since signal propagation at finite width tends to reduce the scale of feature kernels [14]. We suspect that this is the cause of the slower dynamics of ensembled predictors for narrower networks in Figure 7 and Figure A.4. Additionally, the term involving Σ^K will generically increase the cross term since the dynamics of Δ cause its fluctuations to become anti-correlated with the fluctuations in K . In general, it is challenging to make strong definitive statements about the relative scale of these competing effects on the cross term. However, we can say more about this solution in the lazy limit, where we find that the cross term will generically be positive, leading to larger MSE (Appendix H.2).

G.2 Perturbation Theory in Rates rather than Predictions

In experiments on deep CNNs trained on CIFAR-10 in 7 and A.4, we find that the loss curves for the ensemble averaged predictors are effectively time rescaled by a function of network width. In this section, we argue that a proper way to account for this is to compute a perturbation expansion in the *exponent* which defines the rate of decay of the training errors. To illustrate the point, we first consider the case of a single training example before describing larger datasets. In this case, we consider the change of variables $\Delta(t) = e^{-r(t)y}$. We now treat r as an order parameter of the theory with dynamics

$$\frac{d}{dt} r(t) = K(t) \quad (53)$$

Note that this equation is now a linear relation between two order parameters $(r(t), K(t))$, whereas the relation was previously quadratic. In the lazy limit, if $K \rightarrow K - \epsilon$ then $r \rightarrow r - \epsilon t$, giving an effective rescaling of training time by $1 - \bar{\kappa}$.

For multiple training examples, we introduce the notion of a transition matrix $T(t) \in \mathbb{R}^{P \times P}$ which has dynamics

$$\frac{d}{dt} T(t) = -K(t) T(t), \quad T(0) = \mathbf{I}. \quad (54)$$

The solution to the training prediction errors can be obtained at any time t by multiplying the initial condition $\Delta(0) = \mathbf{y}$ with the transition matrix $\Delta(t) = T(t)\mathbf{y}$, where \mathbf{y} are the training targets. In this case, the relevant *rate matrix*, which would be an alternative order parameter is

$$\mathbf{R}(t) = -\log T(t) \quad (55)$$

where \log is the matrix logarithm function. Note that in general $T(t)$ admits a Peano-Baker series solution [62–64]. In the special case where $\mathbf{K}(t)$ commutes with $\bar{\mathbf{K}}(t) = \frac{1}{t} \int_0^t ds \mathbf{K}(s)$, we obtain the following simplified formula for the rate matrix \mathbf{R}

$$\mathbf{R}(t) = \int_0^t ds \mathbf{K}(s) \quad (56)$$

The benefit of this representation is the elimination of coupled order parameter dynamics which are quadratic in fluctuations (in \mathbf{R} and \mathbf{K}) into a linear dynamical relation between order parameters \mathbf{R} and \mathbf{K} . An expansion in \mathbf{R} will thus give better predictions at long times t than a direct expansion in \mathbf{K} . In the lazy $\gamma \rightarrow 0$ limit, the constancy of $\mathbf{K}(t) = \mathbf{K}$ gives the further simplification $\mathbf{R} = \mathbf{K}t$. Working with this representation, we have the following finite width expression for the training loss

$$\begin{aligned} |\langle \mathbf{R}(t) \rangle|^2 &= \mathbf{y}^\top \langle \exp(-2\mathbf{R}(t)) \rangle \mathbf{y} \\ &\sim \mathbf{y}^\top \exp\left[-2\mathbf{R}_1(t) + \frac{1}{N}\mathbf{R}^1(t)\right] \mathbf{y} \\ &+ \frac{1}{2} \times \Sigma^{\mathbf{R}}(\mathbf{R}, t) \frac{\partial^2}{\partial \mathbf{R} \partial \mathbf{R}} \mathbf{y}^\top \exp(-2\mathbf{R}) \mathbf{y} \Big|_{\mathbf{R}=\mathbf{R}_1(t) + \frac{1}{N}\mathbf{R}^1(t)} + \mathcal{O}(N^{-2}) \end{aligned} \quad (57)$$

where $\langle \mathbf{R} \rangle \sim \mathbf{R}_1 + \frac{1}{N}\mathbf{R}^1 + \mathcal{O}(N^{-2})$ is the leading correction to the mean \mathbf{R} . In this representation, it is clear that finite width can alter the timescale of the dynamics through a correction to the mean of \mathbf{R} , as well as contribute an additive correction from fluctuations. This justifies the study perturbation analysis of rates R_N as a function of $1/N$ in Figures 7 and A.4.

H Variance in the Lazy Limit

We can simplify the propagator equations in the lazy $\gamma \rightarrow 0$ limit. To demonstrate how to use our formalism, we go through the complete process of inverting the Hessian, however, for this case, this procedure is a bit cumbersome. A simplified derivation for the lazy limit can be found below in section H.1 which relies only on linearizing the dynamics around the infinite width solution. In the $\gamma \rightarrow 0$ limit, all of the D tensors vanish and the κ tensors are constant in time. Thus, it suffices to analyze the kernels restricted to $t = 0$ and study the evolution of the prediction variance $\Delta(t)$.

$$\begin{aligned} S &= \int dt \times \hat{\Delta}(t) \Delta(t) - \mathbf{y} + \int ds \times \Theta(t-s) \mathbf{K} \Delta(s) \\ &+ \times \times \mathbf{h} \cdot \hat{\Phi} \cdot \Phi + G \cdot \hat{G} \cdot \mathbf{i} + \times \hat{K} \cdot \mathbf{K} - \times G^{+1} \Phi \cdot + \times \ln \mathcal{Z} \\ \mathcal{Z} &= \mathbb{E}_{\mathbf{r} \sim \mathcal{N}(0, \mathbf{G}^{-1}), \mathbf{g} \sim \mathcal{N}(0, \mathbf{G}^{+1})} \exp\left[-\hat{\Phi} \cdot \phi(\mathbf{u}) \phi(\mathbf{u}) - \hat{G} \cdot \mathbf{g} \cdot \mathbf{g}\right], \mathbf{g} = \mathbf{r} \cdot \dot{\phi}(\mathbf{u}) \end{aligned} \quad (58)$$

where $\{\mathbf{u}\} \sim \mathcal{N}(0, \mathbf{I}^{-1})$, $\{\mathbf{r}\} \sim \mathcal{N}(0, \mathbf{G}^{+1})$. Taking two derivatives with respect to $\{\hat{\Phi}, \hat{G}\}$ give terms of the form

$$\begin{aligned} \kappa &= \phi(\mathbf{u}) \phi(\mathbf{u}) \phi(\mathbf{u}) \phi(\mathbf{u}) - \Phi \cdot \Phi \\ \kappa^{\mathbf{G}} &= \mathbf{g} \cdot \mathbf{g} \cdot \mathbf{g} \cdot \mathbf{g} - \mathbf{G} \cdot \mathbf{G} \\ \kappa^{:\mathbf{G}} &= \phi(\mathbf{u}) \phi(\mathbf{u}) \mathbf{g} \cdot \mathbf{g} - \Phi \cdot \mathbf{G} \end{aligned} \quad (59)$$

Given these we also have the relevant non-vanishing sensitivity tensors

$$\begin{aligned}
D^{\hat{u}^i+1} &= \frac{\partial^2}{\partial \Phi^i} \phi(u^i+1)\phi(u^i+1), \quad D^{G^i G^i+1} = \frac{\partial}{\partial G^i+1} g^i g^i \\
D^{G^i \hat{u}^i+1} &= \frac{\partial}{\partial \Phi^i+1} g^i g^i \\
D^{K^i} &= \delta_{ij} \delta^i G^i+1, \quad D^{K^i G^i} = \delta^i \delta^i \Phi^i+1 \\
D^{-K}(t) &= \int ds \Theta(t-s) \delta^i \Delta^i(s)
\end{aligned} \tag{60}$$

As before we let $q_1 = \text{Vec}\{\Delta^i(t), \hat{\Phi}^i, G^i, K^i\}$ and $q_2 = \text{Vec}\{\hat{\Delta}^i(t), \hat{\Phi}^i, \hat{G}^i, \hat{K}^i\}$. The propagator has the form

$$\mathbf{U} \equiv \nabla_{q_2 q_1}^2 S = \begin{pmatrix} \mathbf{I} + \mathbf{K} & \mathbf{0} & \mathbf{0} & D^{-K} \\ \mathbf{0} & \mathbf{I} - D^i & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & -D^G & \mathbf{I} - D^{GG} & \mathbf{0} \\ \mathbf{0} & -D^K & -D^{KG} & \mathbf{I} \end{pmatrix}, \quad \nabla_{q_2 q_2}^2 S = \begin{pmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & G & G & \mathbf{0} \\ \mathbf{0} & G & GG & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix} \tag{62}$$

The propagator of interest is $\langle q_1 = \mathbf{U}^{-1} \nabla_{q_2 q_2}^2 S \mathbf{U}^{-1} \rangle$. We can exploit the block structure of \mathbf{U} to find an inverse

$$\mathbf{U}^{-1} = \begin{pmatrix} \mathbf{U}^{-1} & \mathbf{U}^{-1} & \mathbf{U}_G^{-1} & \mathbf{U}_K^{-1} \\ \mathbf{0} & \mathbf{U}^{-1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{U}_G^{-1} & \mathbf{U}_{GG}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{U}_K^{-1} & \mathbf{U}_{KG}^{-1} & \mathbf{I} \end{pmatrix} \tag{63}$$

where each sub-block can be computed with the Schur-complement formula. Altogether, we multiply through to get the propagator

$$\begin{aligned}
& \begin{pmatrix} \mathbf{0} & \mathbf{U}^{-1} & \mathbf{U}_G^{-1} & \mathbf{U}_K^{-1} \\ \mathbf{0} & \mathbf{U}_G^{-1} & \mathbf{U}_{GG}^{-1} & \mathbf{U}_K^{-1} \\ \mathbf{0} & \mathbf{U}_K^{-1} & \mathbf{U}_{KG}^{-1} & \mathbf{I} \end{pmatrix} + \begin{pmatrix} \mathbf{U}_G^{-1} & \mathbf{U}_{GG}^{-1} & \mathbf{U}_K^{-1} \\ \mathbf{U}_G^{-1} & \mathbf{U}_{GG}^{-1} & \mathbf{U}_K^{-1} \\ \mathbf{U}_K^{-1} & \mathbf{U}_{KG}^{-1} & \mathbf{I} \end{pmatrix} \\
& \times \begin{pmatrix} [\mathbf{U}^{-1}] & [\mathbf{U}_G^{-1}] & [\mathbf{U}_K^{-1}] \\ [\mathbf{U}_G^{-1}] & [\mathbf{U}_{GG}^{-1}] & [\mathbf{U}_{KG}^{-1}] \\ [\mathbf{U}_K^{-1}] & [\mathbf{U}_{KG}^{-1}] & \mathbf{I} \end{pmatrix}
\end{aligned} \tag{64}$$

Two of these blocks corresponding to K, Δ are especially important for characterizing the fluctuations of network predictions. The covariance structure for K has the form

$$\kappa = \mathbf{U}_K^{-1} [\mathbf{U}_K^{-1}] + \mathbf{U}_{KG}^{-1} G [\mathbf{U}_K^{-1}] + \mathbf{U}_K^{-1} G [\mathbf{U}_{KG}^{-1}] + \mathbf{U}_{KG}^{-1} GG [\mathbf{U}_{KG}^{-1}] \tag{65}$$

Next we use the fact that $\mathbf{U}^{-1} = \mathbf{U}_K^{-1} \mathbf{U}_K^{-1}$ and that $\mathbf{U}_G^{-1} = \mathbf{U}_K^{-1} \mathbf{U}_{KG}^{-1}$, which follows from the block structure of \mathbf{U} . Consequently we arrive at the identity

$$\begin{aligned}
& = \mathbf{U}^{-1} [\mathbf{U}^{-1}] + \mathbf{U}_G^{-1} G [\mathbf{U}^{-1}] + \mathbf{U}_G^{-1} G [\mathbf{U}^{-1}] + \mathbf{U}_G^{-1} GG [\mathbf{U}_G^{-1}] \\
& = \mathbf{U}_K^{-1} \kappa [\mathbf{U}_K^{-1}].
\end{aligned} \tag{66}$$

Lastly, we note that, by the Schur-complement formula that $\mathbf{U}_K^{-1} = -(\mathbf{I} + \kappa)^{-1} D^{-K}$. Thus, writing $(\mathbf{I} + \kappa)^{-1} (\mathbf{I} + \kappa)^{-1} = D^{-K} \kappa [D^{-K}]$ as an integral equation, we find

$$\begin{aligned}
\Sigma(t, s) &+ \int_0^t dt^\theta \int_0^s ds^\theta K \Sigma(t^\theta, s) + \int_0^s ds^\theta \int_0^t dt^\theta K \Sigma(t, s^\theta) \\
&+ \int_0^t dt^\theta \int_0^s ds^\theta K K \Sigma(t^\theta, s^\theta) = \int_0^t \Delta(t^\theta) \int_0^s ds^\theta \Delta(s^\theta) \Sigma^K;
\end{aligned} \tag{67}$$

Differentiation with respect to t and s gives a simple differential equation

$$\begin{aligned} \frac{\partial^2}{\partial t \partial s} \Sigma(t, s) + \mathcal{K} \frac{\partial}{\partial s} \Sigma(t, s) + \mathcal{K} \frac{\partial}{\partial t} \Sigma(t, s) \\ + \mathcal{K} \Sigma(t, s) = \Delta(t) \Delta(s) \Sigma^K; \end{aligned} \quad (68)$$

Let $\{\kappa_j\}$ be the eigenvectors of the kernel matrix \mathcal{K} . Projecting these dynamics on the eigenspace $\Sigma_{\kappa^0}(t, s) = \kappa^0(t, s)$ recovers the equation in the main text

$$\left(\frac{\partial}{\partial t} + \lambda_{\kappa^0} \right) \left(\frac{\partial}{\partial s} + \lambda_{\kappa^0} \right) \Sigma_{\kappa^0}(t, s) = \Delta_{\kappa^0}(t) \Delta_{\kappa^0}(s) \Sigma_{\kappa^0}^K. \quad (69)$$

Replacing $\Sigma^K = \kappa$ recovers the equation (7) in the main text.

H.1 Perturbed Linear System

In this section, we provide a simpler derivation of the lazy limit training error variance dynamics. In this case, we merely perturb the dynamics around its infinite width value $\epsilon(t) = \epsilon^1(t) + \epsilon^K(t)$ and $\mathcal{K} = \mathcal{K}^1 + \mathcal{K}^K$, and keep terms only linear in these perturbations. The perturbation ϵ^K is fixed in time and the dynamics of $\epsilon^1(t)$ are

$$\frac{d}{dt} \epsilon^1(t) = -\mathcal{K}^1 \epsilon^1(t) - \mathcal{K}^K \epsilon^1(t) \quad (70)$$

Projecting this equation on the eigenspace of \mathcal{K}^1 gives

$$\frac{d}{dt} \epsilon_{\kappa^0}^1(t) = -\lambda_{\kappa^0} \epsilon_{\kappa^0}^1(t) - \epsilon_{\kappa^0}^K \Delta_{\kappa^0}^1(t) \quad (71)$$

This immediately recovers the final result of the last section

$$\begin{aligned} N \left(\frac{\partial}{\partial t} + \lambda_{\kappa^0} \right) \left(\frac{\partial}{\partial s} + \lambda_{\kappa^0} \right) \epsilon_{\kappa^0}^1(t) \epsilon_{\kappa^0}^1(s) &= \left(\frac{\partial}{\partial t} + \lambda_{\kappa^0} \right) \left(\frac{\partial}{\partial s} + \lambda_{\kappa^0} \right) \Sigma_{\kappa^0}(t, s) \\ &= \Sigma_{\kappa^0}^K \Delta_{\kappa^0}^1(t) \Delta_{\kappa^0}^1(s) \end{aligned} \quad (72)$$

Qualitatively, the process of computing this linear correction (in ϵ^K) to the dynamics of ϵ^1 is identical to the argument utilized in prior work on perturbative feature learning corrections [11]. In that context, the perturbation is caused by small amounts of feature learning, rather than initialization fluctuations.

H.2 Mean Prediction Error Correction in the Lazy Limit

Using a similar heuristic as in the preceding section, we now consider the correction to the mean predictor $\langle \Delta(t) \rangle$ in the lazy limit. Taylor expanding $\langle \Delta(t) \rangle$ in powers of $1/N$, we find

$$\begin{aligned} \frac{d}{dt} \langle \Delta(t) \rangle &= \frac{d}{dt} \epsilon^1(t) + \frac{1}{N} \frac{d}{dt} \epsilon^K(t) + \dots \\ &= -\langle (\mathcal{K} - \mathcal{K}^1 + \mathcal{K}^K) (\epsilon^1 + \epsilon^K) \rangle \\ &= -\mathcal{K}^1 \epsilon^1 - \mathcal{K}^K \langle \epsilon^1 \rangle \\ &\quad - \langle (\mathcal{K} - \mathcal{K}^1) \rangle \epsilon^1 - \langle (\mathcal{K} - \mathcal{K}^1) \rangle (\epsilon^1 + \epsilon^K) \\ &\sim -\mathcal{K}^1 \epsilon^1 - \frac{1}{N} \mathcal{K}^K \epsilon^1 - \frac{1}{N} \mathcal{K}^1 \epsilon^1 - \frac{1}{N} \mathcal{K}^K \epsilon^1 + \mathcal{O}(N^{-2}) \end{aligned} \quad (73)$$

From the previous section we have that

$$\frac{d}{dt} \epsilon^1 = -\mathcal{K}^1 \epsilon^1 - \mathcal{K}^K \epsilon^1 \implies \epsilon^1(t) = - \int_0^t ds \exp(-\mathcal{K}^1(t-s)) \mathcal{K}^K \exp(-\mathcal{K}^1 s) \mathbf{y} \quad (74)$$

Projecting these dynamics onto the eigenspace of the kernel gives

$$\epsilon_k(t) = - \sum_k \epsilon_k^K \frac{e^{-\lambda_k t} - e^{-\lambda_\ell t}}{\lambda_k - \lambda_\ell} y \quad (75)$$

where $\ell = k$ should be seen as the limit where $\lambda_k \rightarrow \lambda_\ell$ of the above. Thus we find that the leading mean correction to the error solves the following differential equation

$$\begin{aligned} \frac{d}{dt} + \lambda_k \Delta_k^1(t) &= - \sum_k K_k^1 y e^{-\lambda_k t} + \sum_{k \neq \ell} \epsilon_k^K \frac{e^{-\lambda_\ell t} - e^{-\lambda_k t}}{\lambda_\ell - \lambda_k} y \\ &= \sum_k y e^{-\lambda_k t} - K_k^1 + \sum_k \epsilon_k^K t + \sum_{k \neq \ell} \epsilon_k^K \frac{e^{-\lambda_\ell t} - e^{-\lambda_k t}}{\lambda_\ell - \lambda_k} y \end{aligned} \quad (76)$$

We see that at late sufficiently large t , that the terms involving Σ^K will dominate. We can gain more intuition by considering the special case of a single training data point where the mean error correction has the form

$$\begin{aligned} \frac{d}{dt} + \lambda \Delta^1(t) &= y e^{-\lambda t} - K^1 + t \Sigma^K \implies \Delta^1(t) = y e^{-\lambda t} - K^1 + \frac{1}{2} t^2 \Sigma^K e^{-\lambda t} \\ \implies \Delta(t)^2 &\sim \Delta^1(t)^2 + \frac{1}{N} 2y^2 t e^{-2\lambda t} - K^1 + \frac{1}{2} t \Sigma^K + \Sigma(t, t) + \mathcal{O}(N^{-2}) \\ &\sim \Delta^1(t)^2 + \frac{2}{N} y^2 t e^{-2\lambda t} - K^1 + \Sigma^K t + \mathcal{O}(N^{-2}) \end{aligned} \quad (77)$$

While the term involving Σ^K is positive for all t , K^1 could be positive or negative for a given architecture. If K^1 is positive, then MSE is initially improved at early times but after $t > \frac{K^1}{\Sigma^K}$ the MSE is worse than the infinite width. On the other hand, if K^1 is negative (as we suspect is typically the case), then the MSE will strictly decrease with network width for any time t .

I Two Layer Equations and Time/Time Diagonal

In this section, we analyze two layer networks in greater detail. Unlike the deep network case, two layer networks can be analyzed on the time-time diagonal: ie the dynamics only depend on $\Phi(t, t)$ and $G(t, t)$ rather than on all possible off-diagonal pairs of time points. Further, there are no response functions A, B which complicate the recipe for calculating the propagator (Appendix E).

I.1 A Single Training Point

For a two layer network trained on a single training point with norm constraint $|\mathbf{x}|^2 = D$, we have the following DMFT action

$$\begin{aligned} S[\{K(t), \hat{K}(t), \Delta(t), \hat{\Delta}(t)\}] &= \int dt [K(t) \hat{K}(t) + \hat{\Delta}(t) \Delta(t) - y + \int ds \Theta(t-s) \Delta(s) K(s)] \\ &+ \ln \mathcal{Z}[\hat{K}, f], \quad \mathcal{Z} = \int \mathcal{D}h, g \exp \left[- \int dt \hat{K}(t) [\phi(h(t))^2 + g(t)^2] \right] \end{aligned} \quad (78)$$

The saddle point equations are

$$\begin{aligned} \frac{\partial S}{\partial \hat{K}(t)} &= K(t) - [\phi(h(t))^2 + g(t)^2] = 0 \\ \frac{\partial S}{\partial \hat{\Delta}(t)} &= \Delta(t) - y + \int ds \Theta(t-s) \Delta(s) K(s) = 0 \\ \frac{\partial S}{\partial K(s)} &= \hat{K}(s) + \Delta(s) \int dt \hat{\Delta}(t) \Theta(t-s) = 0 \\ \frac{\partial S}{\partial \Delta(s)} &= \hat{\Delta}(s) + K(s) \int dt \hat{\Delta}(t) \Theta(t-s) = 0 \end{aligned} \quad (79)$$

From these equations, we can compute the entries in the Hessian of the DMFT action S . Letting

$$\mathbf{q}(t) = \begin{pmatrix} \Delta(t) \\ K(t) \end{pmatrix} \text{ and } \hat{\mathbf{q}}(t) = \begin{pmatrix} \hat{\Delta}(t) \\ \hat{K}(t) \end{pmatrix}$$

$$\begin{aligned} \frac{\partial^2 S}{\partial \mathbf{q}(t) \partial \mathbf{q}(s)^T} &= \mathbf{0} \\ \frac{\partial^2 S}{\partial \hat{\mathbf{q}}(t) \partial \mathbf{q}(s)^T} &= \begin{pmatrix} \delta(t-s) + \Theta(t-s)K(s) & \Theta(t-s)\Delta(s) \\ -\frac{\partial}{\partial (s)}(\phi(h(t))^2 + g(t)^2) & \delta(t-s) \end{pmatrix} \\ \frac{\partial^2 S}{\partial \hat{\mathbf{q}}(t) \partial \hat{\mathbf{q}}(s)^T} &= \begin{pmatrix} 0 & 0 \\ 0 & \kappa(t,s) \end{pmatrix} \end{aligned} \quad (80)$$

where $\kappa(t,s) = (\phi(h(t))^2 + g(t)^2)(\phi(h(s))^2 + g(s)^2) - K(t)K(s)$ is the NTK's fourth cumulant. We now vectorize our order parameters over time $\mathbf{q} = \text{Vec}\{\mathbf{q}(t)\}_{t \in \mathbb{R}_+}$ and $\hat{\mathbf{q}} = \text{Vec}\{\hat{\mathbf{q}}(t)\}_{t \in \mathbb{R}_+}$ and express the full Hessian

$$\nabla^2 S = \begin{pmatrix} \mathbf{0} & \frac{\partial^2 S}{\partial \mathbf{q} \partial \hat{\mathbf{q}}^T} \\ \frac{\partial^2 S}{\partial \hat{\mathbf{q}} \partial \mathbf{q}^T} & \frac{\partial^2 S}{\partial \hat{\mathbf{q}} \partial \hat{\mathbf{q}}^T} \end{pmatrix} \Rightarrow -[\nabla^2 S]^{-1} = \begin{pmatrix} (\frac{\partial^2 S}{\partial \hat{\mathbf{q}} \partial \mathbf{q}^T})^{-1} & \frac{\partial^2 S}{\partial \hat{\mathbf{q}} \partial \mathbf{q}^T} (\frac{\partial^2 S}{\partial \mathbf{q} \partial \hat{\mathbf{q}}^T})^{-1} & -(\frac{\partial^2 S}{\partial \hat{\mathbf{q}} \partial \mathbf{q}^T})^{-1} \\ -(\frac{\partial^2 S}{\partial \mathbf{q} \partial \hat{\mathbf{q}}^T})^{-1} & \mathbf{0} & \mathbf{0} \end{pmatrix} \quad (81)$$

The covariance matrix of interest (for $\mathbf{q}(t)$) is thus

$$\mathbf{q} = \begin{pmatrix} \mathbf{I} + \mathbf{K} & \mathbf{0} \\ -\mathbf{D} & \mathbf{I} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mathbf{I} + \mathbf{K} & \mathbf{0} \\ -\mathbf{D} & \mathbf{I} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix} \quad (82)$$

where $[\mathbf{K}](t,s) = \Theta(t-s)K(s)$ and $[\mathbf{D}](t,s) = \Theta(t-s)\Delta(s)$. The above equations allow one to use the infinite width DMFT dynamics for $K(t), \Delta(t)$ to compute the finite size fluctuation dynamics of the kernel K and the error signal Δ .

I.1.1 Computing Field Sensitivities

In this section, we compute $D(t,s)$ by solving for the sensitivity of order parameters. We start with the DMFT field equations

$$h(t) = u + \gamma \int_0^t ds \Delta(s) g(s), \quad z(t) = r + \gamma \int_0^t ds \Delta(s) \phi(h(t)). \quad (83)$$

Now, differentiating both sides with respect to $\Delta(s^\theta)$ gives

$$\begin{aligned} \frac{\partial h(t)}{\partial \Delta(s^\theta)} &= \gamma \Theta(t-s^\theta) g(s^\theta) + \gamma \int_0^t ds \Delta(s) \frac{\partial g(s)}{\partial \Delta(s^\theta)} \\ \frac{\partial z(t)}{\partial \Delta(s^\theta)} &= \gamma \Theta(t-s^\theta) \phi(h(s^\theta)) + \gamma \int_0^t ds \Delta(s) \frac{\partial \phi(h(s))}{\partial \Delta(s^\theta)}. \end{aligned} \quad (84)$$

We can compute D Monte carlo by iteratively solving the above equations for each sampled trajectory $\{h(t), z(t)\}$ [65, 46]. Averaging the necessary fields over the Monte Carlo samples will give us the final expressions for $D(t,s)$.

$$D(t,s) = \frac{\partial}{\partial \Delta(s)} (\phi(h(t))^2 + g(t)^2) \quad (85)$$

Similarly, the uncoupled kernel variance $\kappa(t,s)$ can be evaluated via Monte Carlo sampling for nonlinear networks.

I.2 Test Point Fluctuation Dynamics

We now are in a position to calculate the test/train kernel and test prediction fluctuations. To do this systematically, we augment S with the test point prediction f_γ and field h_γ and introduce the kernel

$K_{\gamma}(t) = \langle \phi(h(t))\phi(h_{\gamma}(t)) + g(t)g_{\gamma}(t) \rangle$. The test prediction f_{γ} and field h_{γ} have dynamics

$$\begin{aligned} h_{\gamma}(t) &= u_{\gamma} + \gamma \int_0^t ds \Delta(s) \dot{\phi}(h_{\gamma}(s)) z(s) K_{\gamma}^x, \quad \langle u_{\gamma} u \rangle = K_{\gamma}^x \\ \frac{\partial}{\partial t} f_{\gamma}(t) &= K_{\gamma}(t) \Delta(t), \quad K_{\gamma}(t) = \langle \phi(h(t))\phi(h_{\gamma}(t)) + g(t)g_{\gamma}(t) \rangle \end{aligned} \quad (86)$$

The augmented action for this DMFT has the form

$$\begin{aligned} S &= \int dt \hat{f}_{\gamma}(t) f_{\gamma}(t) - \int ds \Theta(t-s) \Delta(s) K_{\gamma}(s) + \int dt \hat{K}_{\gamma}(t) K_{\gamma}(t) \\ &+ \int dt \hat{\Delta}(t) \Delta(t) - y + \int ds \Theta(t-s) \Delta(s) K(s) + \int dt \hat{K}(t) K(t) \\ &+ \ln E \exp - \int \hat{K}(t) (\phi(h(t))^2 + g(t)^2) - \int \hat{K}_{\gamma}(t) (\phi(h(t))\phi(h_{\gamma}(t)) + g(t)g_{\gamma}(t)) \end{aligned} \quad (87)$$

We let $\mathbf{q}(t) = [\Delta(t), f_{\gamma}(t), K(t), K_{\gamma}(t)]^>$

$$\begin{aligned} \nabla_{\mathbf{q}\mathbf{q}}^2 S[\mathbf{q}, \hat{\mathbf{q}}] &= \begin{pmatrix} 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & ? & ? \\ 0 & 0 & ? & ?? \end{pmatrix}, \quad \nabla_{\hat{\mathbf{q}}\hat{\mathbf{q}}}^2 S[\mathbf{q}, \hat{\mathbf{q}}] = \begin{pmatrix} 2 & \mathbf{1} + \kappa & 0 & 0 \\ 0 & -\mathbf{D}_{\gamma} & \mathbf{1} & 0 \\ 0 & 0 & \mathbf{1} & 0 \\ 0 & -\mathbf{D}_{\gamma} & 0 & \mathbf{1} \end{pmatrix} \end{aligned} \quad (88)$$

$$D(t, s) = \frac{\partial}{\partial \Delta(s)} (\phi(h(t))^2 + g(t)^2) \quad (88)$$

$$D_{\gamma}(t, s) = \frac{\partial}{\partial \Delta(s)} (\phi(h(t))\phi(h_{\gamma}(t)) + g(t)g_{\gamma}(t)) \quad (89)$$

Our total covariance matrix / propagator is thus

$$\begin{aligned} &= \begin{pmatrix} 2 & \mathbf{1} + \kappa & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & -\mathbf{D}_{\gamma} & \mathbf{1} & 0 & - & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \mathbf{1} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \mathbf{1} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \mathbf{1} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \mathbf{1} & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \mathbf{1} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \mathbf{1} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \mathbf{1} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \mathbf{1} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \mathbf{1} \end{pmatrix} \end{aligned} \quad (90)$$

This is the equation provided in the main text Equation (8).

I.3 Two Layer Linear Network Closed Form

For a linear network on a single data point, we can compute $D(t, s)$ and $\kappa(t, s)$ analytically. We start from the field equations

$$\frac{dh(t)}{dt} = \gamma \Delta(t) z(t), \quad \frac{dz(t)}{dt} = \gamma \Delta(t) h(t) \quad (91)$$

We can make a change of variables $v_+(t) = \frac{1}{\sqrt{2}}(h(t) + z(t))$ and $v_-(t) = \frac{1}{\sqrt{2}}(h(t) - z(t))$. We note that $v_+(0) = \frac{1}{\sqrt{2}}(u + r)$ and $v_-(0) = \frac{1}{\sqrt{2}}(u - r)$ are independent Gaussians. These functions $v_+(t), v_-(t)$ satisfy dynamics

$$\begin{aligned} \frac{dv_+}{dt} &= \gamma \Delta(t) v_+(t), \quad \frac{dv_-}{dt} = -\gamma \Delta(t) v_-(t) \\ \implies v_+(t) &= \exp \left(\gamma \int_0^t ds \Delta(s) \right) v_+(0) \implies \frac{\partial v_+(t)}{\partial \Delta(s)} = \gamma v_+(t) \Theta(t-s) \\ \implies v_-(t) &= \exp \left(-\gamma \int_0^t ds \Delta(s) \right) v_-(0) \implies \frac{\partial v_-(t)}{\partial \Delta(s)} = -\gamma v_-(t) \Theta(t-s) \end{aligned} \quad (92)$$

Now, we use the fact that $v_+(0) = \frac{1}{\sqrt{2}}(u+r)$ and $v_-(0) = \frac{1}{\sqrt{2}}(u-r)$ are independent standard normal random variables to compute $K(t) = \langle h(t)^2 + z(t)^2 \rangle = \langle v_+(t)^2 + v_-(t)^2 \rangle$

$$\begin{aligned} D(t,s) &= \frac{\partial}{\partial \Delta(s)} \langle h(t)^2 + z(t)^2 \rangle = 2\gamma \int_0^t v_+(t)^2 - v_-(t)^2 \Theta(t-s) \\ &= 2\gamma \exp\left(-2\gamma \int_0^t ds \Delta(s)\right) - \exp\left(-2\gamma \int_0^t ds \Delta(s)\right) \Theta(t-s) \end{aligned} \quad (93)$$

This operator is causal ($D(t,s) = 0$ for $s > t$) as expected and vanishes as $t \rightarrow 0$. If we take $\gamma \rightarrow 0$, we have $D(t,s) \rightarrow 0$ which agrees with our reasoning that fields h, z only depend on Δ in the feature learning regime. Since all fields are Gaussian in the linear network case, we can use Wick's theorem to obtain the exact uncoupled kernel variance in the two layer case.

$$\begin{aligned} \kappa(t,s) &= \langle (h(t)^2 + z(t)^2)(h(s)^2 + z(s)^2) \rangle - K(t)K(s) \\ &= 2 \langle h(t)h(s) \rangle^2 + 2 \langle h(t)z(s) \rangle^2 + 2 \langle z(t)h(s) \rangle^2 + 2 \langle z(t)z(s) \rangle^2 \\ &= \langle v_+(t)v_+(s) + v_-(t)v_-(s) \rangle^2 + \langle v_+(t)v_-(s) - v_-(t)v_+(s) \rangle^2 \end{aligned} \quad (94)$$

The $v_\pm(t)$ functions are those given above. Using the fact that $\langle v_+(0)^2 \rangle = \langle v_-(0)^2 \rangle = 1$ allows us to easily compute the single site average above.

J Multiple Samples with Whitened Data

In this section, we analyze the role that sample number plays in dynamics in a simplified model of a two layer linear network trained on whitened data. Concretely, we assume that $\frac{\mathcal{X} \mathcal{X}^T}{D} = \delta$. The field equations for preactivations $h(t)$ and pregradients $z(t)$ obey

$$\frac{d}{dt} h(t) = \gamma \Delta(t) z(t), \quad \frac{d}{dt} z(t) = \gamma \int_0^t \Delta(t) h(s) ds \quad (95)$$

We will assume the targets have unit norm $\|\mathbf{y}\|^2 = 1$ and we define the projection of \mathbf{y} onto the target as $\Delta_y(t) = \mathbf{y} \cdot \mathbf{e}(t)$. The other $P-1$ orthogonal components are denoted $\mathbf{e}_\gamma(t)$ so that $\mathbf{e}(t) = \Delta_y(t) \mathbf{y} + \sum_\gamma \mathbf{e}_\gamma(t)$ with $\mathbf{e}_\gamma(t) \cdot \mathbf{y} = 0$. At infinite width, $\mathbf{e}_\gamma = 0$ and our field equations become

$$\frac{d}{dt} h_y(t) = \Delta_y(t) z(t), \quad \frac{d}{dt} z(t) = \Delta_y(t) h_y(t), \quad \mathbf{e}_\gamma(t) = 0, \quad h_\gamma \sim \mathcal{N}(0,1) \quad (96)$$

However, at finite width N , the off-target predictions \mathbf{e}_γ fluctuate over random initialization. To model all of the fluctuations simultaneously, we consider the following action

$$S = \gamma \int_0^t \int_0^t \hat{\Delta}(t) (\Delta(t) - y) + \ln \mathbb{E} \exp \left[\int_0^t \int_0^t \hat{\Delta}(t) z(t) h(s) \right] \quad (97)$$

which enforces the constraint that $\Delta(t) = y - \frac{1}{N} \langle z(t) h(t) \rangle$ at infinite width. The Hessian over order parameters $\mathbf{q} = \text{Vec}\{\Delta(t), \hat{\Delta}(t)\}$ has the form

$$\nabla_{\mathbf{q}}^2 S = \begin{pmatrix} \mathbf{0} & \\ & (\gamma \mathbf{I} + \mathbf{D}) \end{pmatrix}, \quad D(t,s) = \frac{\partial}{\partial \Delta(s)} z(t) h(t) \quad (98)$$

We thus get the following covariance for predictions $\mathbf{q} = (\gamma \mathbf{I} + \mathbf{D})^{-1} (\gamma \mathbf{I} + \mathbf{D})^{-1}$. We now compute the necessary components of the D tensor

$$\begin{aligned} \frac{\partial h(t)}{\partial \Delta(s)} &= \gamma \delta \Theta(t-s) z(s) + \gamma \int_0^t dt^\theta \Delta(t^\theta) \frac{\partial z(t^\theta)}{\partial \Delta(s)} \\ \frac{\partial z(t)}{\partial \Delta(s)} &= \gamma \Theta(t-s) h(s) + \gamma \int_0^t dt^\theta \Delta(t^\theta) \frac{\partial h(t^\theta)}{\partial \Delta(s)} \\ &= \gamma \Theta(t-s) h(s) + \gamma \int_0^t dt^\theta \Delta_y(t^\theta) \frac{\partial h_y(t^\theta)}{\partial \Delta(s)} \end{aligned} \quad (99)$$

In the last line, we used the fact that these equations are to be evaluated at the mean field infinite width stochastic process where $\Delta_\gamma(t) = 0$. To compute the sensitivity tensor D , we find the following equations for our correlators of interest:

$$\begin{aligned}
\frac{\partial h}{\partial \Delta}(s) z(t) &= \delta_{\mu, \nu} \gamma \Theta(t-s) \langle z(s) z(t) \rangle, \quad \mu, \nu \neq y \\
\frac{\partial z(t)}{\partial \Delta}(s) h(t) &= \gamma \Theta(t-s) \delta_{\mu, \nu} \neq y \\
\frac{\partial h_y(t)}{\partial \Delta_y(s)} z(t) &= \gamma \Theta(t-s) \langle z(s) z(t) \rangle + \gamma \int_0^t dt^\theta \Delta_y(t^\theta) \frac{\partial z(t^\theta)}{\partial \Delta_y(s)} z(t) \\
\frac{\partial z(t)}{\partial \Delta_y(s)} z(t^\theta) &= \gamma \Theta(t-s) \langle h_y(s) z(t) \rangle + \gamma \int_0^t dt^{\theta\theta} \Delta_y(t^{\theta\theta}) \frac{\partial h_y(t^{\theta\theta})}{\partial \Delta_y(s)} z(t^\theta)
\end{aligned} \tag{100}$$

We therefore see that the components of D decouple over indices. In the y direction, we have the following equations

$$D_y(t, s) = \frac{\partial h_y(t)}{\partial \Delta_y(s)} z(t) + \frac{\partial z(t)}{\partial \Delta_y(s)} h_y(t) \tag{101}$$

where the correlators must be solved self-consistently. We will provide this solution in one moment, but first, we will look at the orthogonal directions. For the $P - 1$ orthogonal directions, we obtain the explicit formula for D in each of these directions

$$\begin{aligned}
D_\gamma(t, s) &= \frac{\partial h_\gamma(t)}{\partial \Delta_\gamma(s)} z(t) + \frac{\partial z(t)}{\partial \Delta_\gamma(s)} h_\gamma(t) \\
&= \gamma \Theta(t-s) \langle z(t) z(s) \rangle + \gamma \Theta(t-s)
\end{aligned} \tag{102}$$

Now, we return to D_y . To solve these equations we utilize the change of variables employed in the single sample case $v_+(t) = \frac{1}{\sqrt{2}}(h_y(t) + z(t))$, $v_-(t) = \frac{1}{\sqrt{2}}(h_y(t) - z(t))$ (see Appendix I.3). This orthogonal transformation decouples the dynamics

$$\frac{d}{dt} v_+(t) = \gamma \Delta_y(t) v_+(t), \quad \frac{d}{dt} v_-(t) = -\gamma \Delta_y(t) v_-(t) \tag{103}$$

As a consequence, the field derivatives close

$$\begin{aligned}
\frac{\partial v_+(t)}{\partial \Delta_y(s)} &= \gamma \Theta(t-s) v_+(s) + \int_0^t dt^\theta \Delta_y(t^\theta) \frac{\partial v_+(t^\theta)}{\partial \Delta_y(s)} \\
\frac{\partial v_-(t)}{\partial \Delta_y(s)} &= -\gamma \Theta(t-s) v_-(s) - \int_0^t dt^\theta \Delta_y(t^\theta) \frac{\partial v_-(t^\theta)}{\partial \Delta_y(s)}
\end{aligned} \tag{104}$$

The correlator of interest is

$$\langle h_y(t) z(t) \rangle = \frac{1}{2} \langle [v_+(t) + v_-(t)][v_+(t) - v_-(t)] \rangle = \frac{1}{2} \langle v_+(t)^2 - v_-(t)^2 \rangle \tag{105}$$

So we get that

$$\begin{aligned}
D_y(t, s) &= \frac{1}{2} \frac{\partial}{\partial \Delta_y(s)} \langle v_+(t)^2 - v_-(t)^2 \rangle \\
&= \frac{1}{2} \langle v_+(t) \frac{\partial v_+(t)}{\partial \Delta_y(s)} - v_-(t) \frac{\partial v_-(t)}{\partial \Delta_y(s)} \rangle
\end{aligned} \tag{106}$$

Similarly, we can derive the on-target and off-target uncoupled variances $\kappa_y(t, s)$ and $\kappa_\gamma(t, s)$, which satisfy

$$\begin{aligned}
\kappa_y(t, s) &= \langle v_+(t) v_+(s) + v_-(t) v_-(s) \rangle^2 + \langle v_+(t) v_+(s) - v_-(t) v_-(s) \rangle^2 \\
\kappa_\gamma(t, s) &= \frac{1}{2} \langle v_+(t) v_+(s) + v_-(t) v_-(s) \rangle
\end{aligned} \tag{107}$$

Using these functions, we arrive at the following variance for each of the P dimensions

$$\begin{aligned} \sigma_y &= (\gamma \mathbf{I} + \mathbf{D}_y)^{-1} \sigma_y (\gamma \mathbf{I} + \mathbf{D}_y)^{-1} \\ \sigma_\gamma &= (\gamma \mathbf{I} + \mathbf{D}_\gamma)^{-1} \sigma_\gamma (\gamma \mathbf{I} + \mathbf{D}_\gamma)^{-1} \end{aligned} \quad (108)$$

Using the fact that all Δ_γ variables are independent and identically distributed under the leading order picture, the expected training loss has the form

$$|\sigma|^2 \approx \Delta_y^\top(t)^2 + \frac{2}{N} \Delta_y^\top(t) \Delta_y^\top(t) + \frac{1}{N} \Sigma_y(t, t) + \frac{(P-1)}{N} \Sigma_\gamma(t, t) + \mathcal{O}(N^{-2}). \quad (109)$$

where $\Delta_y - \Delta_y^\top = \frac{1}{N} \Delta_y^\top(t) + \mathcal{O}(N^{-2})$. We note that the bias correction is $\mathcal{O}(N^{-1})$ while the variance is $\mathcal{O}(P/N)$. We compare the above leading order theory with and without the bias correction in Appendix Figure A.2.

K Online Learning

Our technology for computing finite size effects can easily be translated to a setting where the neural network is trained in an online fashion, disregarding the effect of SGD noise. At each step, we compute the gradient over the full data distribution $p(\mathbf{x})$. Focusing on MSE loss, we study the following equation

$$\frac{d}{dt} \Delta(\mathbf{x}, t) = -\mathbb{E}_{\mathbf{x}^0} p(\mathbf{x}^0) K(\mathbf{x}, \mathbf{x}^0; t) \Delta(\mathbf{x}^0, t) \quad (110)$$

where $K(\mathbf{x}, \mathbf{x}^0; t)$ is the dynamic NTK and $\Delta(\mathbf{x}, t) = y(\mathbf{x}) - f(\mathbf{x}, t)$ is the prediction error. In general the distribution involves integration over an uncountable set of possible inputs \mathbf{x} . To remedy this, we utilize a countable orthonormal basis of functions for the data distribution $\{\psi_k(\mathbf{x})\}_{k=1}^J$. For example, if $p(\mathbf{x})$ were the isotropic Gaussian density for $\mathcal{N}(0, \mathbf{I})$, then ψ_k could be Hermite polynomials. We expand Δ and K in this basis ψ_k , and arrive at the following differential equation

$$\frac{d}{dt} \Delta_k(t) = - \sum_{k^0} K_{kk^0}(t) \Delta_{k^0}(t) \quad (111)$$

By orthonormality, the average turned into a sum over all possible orthonormal functions $\{\psi_k\}$. We note that since K is evolving in time, there is not generally a fixed basis of functions that diagonalize K , resulting in the couplings across eigenmodes in Equation (111). Since, in online learning, there is no distinction between the training and test distribution, our error of interest is simply $\mathcal{L}(t) = \sum_k \Delta_k(t)^2$. To obtain the finite size corrections to this quantity, we compute the joint propagator for all variables $\{K_{kk^0}(t), \Delta_{k^0}(t)\}$. If we wanted to pursue a perturbation theory in rates (Appendix G.2), we could again define a transition matrix T and rate matrix $\mathbf{R}(t)$ as

$$\mathbf{R}(t) = -\log T(t), \quad \frac{d}{dt} T_{kk^0}(t) = - \sum_{k^0} K_{kk^0}(t) T_{k^0}(t), \quad T_{k^0}(0) = \delta_{k^0} \quad (112)$$

We can then obtain $\Delta = \exp(-\mathbf{R}(t)) \mathbf{y}$, where $y_k = \mathbb{E}_{\mathbf{x}} \psi_k(\mathbf{x}) y(\mathbf{x})$. Since \mathbf{R} has a finite size mean correction and finite size fluctuations, so too does the error $\Delta_k(t)$ and the loss \mathcal{L} (Appendix G.2).

K.1 Two Layer Networks

In the two layer case, instead of tracking kernels, we could instead deal with the distribution over read-in vectors $\mathbf{w} \in \mathbb{R}^D$ and readout scalars $a \in \mathbb{R}$ as in the original works on mean field networks [6, 66]. When training on the population risk equations for $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I})$

$$\begin{aligned} \frac{d}{dt} \mathbf{w} &= a \mathbb{E}_{\mathbf{x}} \Delta(\mathbf{x}) \dot{\phi}(\mathbf{w} \cdot \mathbf{x}) \mathbf{x} = \mathbb{E}_{\mathbf{x}} \frac{\partial \Delta(\mathbf{x})}{\partial \mathbf{x}} \dot{\phi}(\mathbf{w} \cdot \mathbf{x}) + \mathbb{E} \Delta(\mathbf{x}) \ddot{\phi}(\mathbf{w} \cdot \mathbf{x}) \mathbf{w} \\ \frac{d}{dt} a &= \mathbb{E}_{\mathbf{x}} \Delta(\mathbf{x}) \phi(\mathbf{w} \cdot \mathbf{x}) \end{aligned} \quad (113)$$

The action Z has the form

$$S = \gamma \int dt d\mathbf{x} \hat{\Delta}(t, \mathbf{x}) (\Delta(t, \mathbf{x}) - y(\mathbf{x})) + \ln \mathbb{E}_{a, \mathbf{w}} \exp \int dt d\mathbf{x} \hat{\Delta}(t, \mathbf{x}) a(t) \phi(\mathbf{w}(t) \cdot \mathbf{x}) \quad (114)$$

The Hessian over $\mathbf{q} = \{\Delta(t), \hat{\Delta}(t)\}$ is

$$\nabla^2 S = \begin{pmatrix} 0 & \mathbf{I} + \mathbf{D} \\ \mathbf{I} + \mathbf{D} & \mathbf{D} \end{pmatrix}. \quad (115)$$

where $D(t, \mathbf{x}; s, \mathbf{x}^\ell) = \frac{\partial}{\partial \Delta(s, \mathbf{x})} a(t) \phi(\mathbf{w}(t) \cdot \mathbf{x})$. We can use the following implicit rule

$$\begin{aligned} \frac{\partial a(t)}{\partial \Delta(s, \mathbf{x})} &= \gamma \Theta(t-s) p(\mathbf{x}) \phi(\mathbf{w}(s) \cdot \mathbf{x}) + \gamma \mathbb{E}_{\mathbf{x}^\ell} \int_0^t dt^\ell \Delta(t^\ell, \mathbf{x}^\ell) \dot{\phi}(\mathbf{w} \cdot \mathbf{x}^\ell) \mathbf{x}^\ell \cdot \frac{\partial \mathbf{w}(t)}{\partial \Delta(s, \mathbf{x})} \\ \frac{\partial \mathbf{w}(t)}{\partial \Delta(s, \mathbf{x})} &= \gamma \Theta(t-s) p(\mathbf{x}) a(s) \dot{\phi}(\mathbf{w}(s) \cdot \mathbf{x}) \mathbf{x} \\ &\quad + \gamma \mathbb{E}_{\mathbf{x}^\ell} \int_0^t dt^\ell \Delta(t^\ell, \mathbf{x}^\ell) \frac{\partial a(t^\ell)}{\partial \Delta(s, \mathbf{x})} \dot{\phi}(\mathbf{w} \cdot \mathbf{x}^\ell) + a(t^\ell) \ddot{\phi}(\mathbf{w} \cdot \mathbf{x}^\ell) \frac{\partial \mathbf{w}(t^\ell)}{\partial \Delta(s, \mathbf{x})} \cdot \mathbf{x}^\ell \end{aligned} \quad (116)$$

The above equations could be solved and then used to compute $D(t, \mathbf{x}; s, \mathbf{x}^\ell)$ which must then be inverted to get the observed prediction variance.

K.2 Linear Activations

Using the ideas in the preceding sections, we can make more progress in the case of a two layer linear network in the online learning setting. The key idea is to track the kernel and prediction error projections onto the space of linear functions. In this case we get the following DMFT over the order parameter $\langle \mathbf{w}(t) \mathbf{w}^\top \rangle = \frac{1}{N} \mathbf{W}^\top \mathbf{a} \in \mathbb{R}^D$.

$$\begin{aligned} \frac{d}{dt} a(t) &= \gamma (\langle \mathbf{w}(t) \mathbf{w}^\top \rangle - \langle \mathbf{w}(t) \mathbf{w}^\top \rangle) \cdot \mathbf{w}(t) \\ \frac{d}{dt} \mathbf{w}(t) &= \gamma a(t) (\langle \mathbf{w}(t) \mathbf{w}^\top \rangle - \langle \mathbf{w}(t) \mathbf{w}^\top \rangle) \\ \langle \mathbf{w}(t) \mathbf{w}^\top \rangle &= \frac{1}{\gamma} \langle a(t) \mathbf{w}(t) \rangle \end{aligned} \quad (117)$$

At infinite width, we see that the dynamics can be reduced to tracking the projection of the weights \mathbf{w} and \mathbf{a} on the $\langle \mathbf{w}(t) \mathbf{w}^\top \rangle$ direction. The $D-1$ off-target dimensions vanish $\langle \mathbf{w}(t) \mathbf{w}^\top \rangle = 0$. At infinite width, we arrive at the alignment dynamics studied in prior work [64, 9]

$$\begin{aligned} \frac{d}{dt} \langle \mathbf{w}(t) \mathbf{w}^\top \rangle &= \mathbf{M}(t) (\langle \mathbf{w}(t) \mathbf{w}^\top \rangle - \langle \mathbf{w}(t) \mathbf{w}^\top \rangle) \\ \frac{d}{dt} \mathbf{M}(t) &= \gamma^2 \langle \mathbf{w}(t) \mathbf{w}^\top \rangle (\langle \mathbf{w}(t) \mathbf{w}^\top \rangle)^\top + \gamma^2 \langle \mathbf{w}(t) \mathbf{w}^\top \rangle (\langle \mathbf{w}(t) \mathbf{w}^\top \rangle)^\top \\ &\quad + 2\gamma^2 \langle \mathbf{w}(t) \mathbf{w}^\top \rangle \cdot \langle \mathbf{w}(t) \mathbf{w}^\top \rangle \mathbf{I} \end{aligned} \quad (118)$$

We note that $\langle \mathbf{w}(t) \mathbf{w}^\top \rangle = \beta(t) \langle \mathbf{w}(t) \mathbf{w}^\top \rangle$ and that \mathbf{M} has only one special eigenvector $\langle \mathbf{w}(t) \mathbf{w}^\top \rangle$ with eigenvalue $m_\gamma(t)$. It thus suffices to track evolution in this single direction

$$\frac{d}{dt} \beta(t) = m_\gamma(t) (\beta_\gamma - \beta(t)), \quad \frac{d}{dt} m_\gamma(t) = 4\gamma^2 \beta(t) (\beta_\gamma - \beta(t)) \quad (119)$$

We note that this equation is identical to the differential equation for a single training example in Appendix J. Here $\beta_\gamma - \beta(t)$ plays the role of $\Delta_y(t)$ and $m_\gamma(t)$ plays the role of the kernel $K_y(t)$. A key observation is the conservation law $4\gamma^2 \frac{d}{dt} \beta(t)^2 = \frac{d}{dt} m_\gamma(t)^2$, from which it follows that $m_\gamma(t)^2 - 4 = 4\gamma^2 \beta(t)$ [9]

$$\frac{d}{dt} \beta(t) = 2 \sqrt{1 + \gamma^2 \beta(t)^2} (\beta_\gamma - \beta(t)) \quad (120)$$

This is identical to the differential equations for a single sample (producing prediction $f(t)$ and kernel $K(t)$) if the following substitutions are made

$$f(t) \leftrightarrow \beta(t), \quad K(t) \leftrightarrow m_\gamma(t) \quad (121)$$

We now proceed to compute finite size corrections starting from the action

$$S = \gamma \int dt \hat{w}(t) \cdot \mathbf{w}(t) + \ln \mathbb{E} \exp \left[- \int dt \hat{w}(t) \cdot \mathbf{w}(t) a(t) \right] \quad (122)$$

The necessary ingredients are

$$\begin{aligned} \langle \mathbf{w}(t), \mathbf{w}(s) \rangle &= \langle a(t)a(s) \mathbf{w}(t) \mathbf{w}(s) \rangle - \gamma^2 \langle \mathbf{w}(t), \mathbf{w}(s) \rangle \\ &= \langle a(t)a(s) \rangle \langle \mathbf{w}(t) \mathbf{w}(s) \rangle + \langle a(s) \mathbf{w}(t) \rangle \langle a(t) \mathbf{w}(s) \rangle \in \mathbb{R}^{D \times D} \end{aligned} \quad (123)$$

Similarly we have to compute the sensitivity tensor

$$\mathbf{D}(t, s) = \frac{\partial}{\partial \langle \mathbf{w}(s) \rangle} a(t) \mathbf{w}(t) \in \mathbb{R}^{D \times D} \quad (124)$$

We start from the dynamics

$$\frac{d}{dt} \mathbf{w}(t) = \gamma a(t) (\boldsymbol{\gamma} - \mathbf{w}(t)), \quad \frac{d}{dt} a(t) = \gamma (\boldsymbol{\gamma} - \mathbf{w}(t)) \cdot \mathbf{w}(t) \quad (125)$$

Next, we have to calculate causal derivatives for fields

$$\begin{aligned} \frac{\partial}{\partial \langle \mathbf{w}(s) \rangle} \mathbf{w}(t) &= -\gamma \Theta(t-s) a(s) \mathbf{I} + \gamma \int_0^t dt^\theta (\boldsymbol{\gamma} - \mathbf{w}(t^\theta)) \frac{\partial a(t^\theta)}{\partial \langle \mathbf{w}(s) \rangle} \\ \frac{\partial}{\partial \langle \mathbf{w}(s) \rangle} a(t) &= -\gamma \Theta(t-s) \mathbf{w}(s) + \gamma \int_0^t dt^\theta (\boldsymbol{\gamma} - \mathbf{w}(t^\theta)) \cdot \frac{\partial \mathbf{w}(t^\theta)}{\partial \langle \mathbf{w}(s) \rangle} \end{aligned} \quad (126)$$

Following an identical argument as in J, we see that \mathbf{D} has block diagonal structure with $D_\gamma(t, s)$ on the γ direction and $D_\beta(t, s)$ in any of the $D-1$ remaining directions

$$D_\gamma(t, s) = \frac{\partial}{\partial \beta(s)} a(t) w_\gamma(t), \quad D_\beta(t, s) = \frac{\partial}{\partial \beta_\beta(s)} a(t) w_\beta(t) \quad (127)$$

Similarly, $\langle \mathbf{w}(t), \mathbf{w}(s) \rangle$ has a similar decomposition

$$\begin{aligned} \kappa_\gamma(t, s) &= \langle a(t)a(s) \rangle \langle w_\gamma(t) w_\gamma(s) \rangle + \langle a(s) w_\beta(t) \rangle \langle a(t) w_\beta(s) \rangle \\ \kappa_\beta(t, s) &= \langle a(t)a(s) \rangle \langle w_\beta(t) w_\beta(s) \rangle + \langle a(s) w_\gamma(t) \rangle \langle a(t) w_\gamma(s) \rangle \end{aligned} \quad (128)$$

The processes have the following equations at infinite width

$$\frac{d}{dt} w_\gamma(t) = \gamma a(t) (\beta_\gamma - \beta(t)), \quad \frac{d}{dt} a(t) = \gamma w_\beta(t) (\beta_\gamma - \beta(t)), \quad \frac{d}{dt} w_\beta(t) = 0 \quad (129)$$

As a consequence we note that $\langle w_\beta(t) a(s) \rangle = 0$ so that $\kappa_\beta(t, s) = \langle a(t)a(s) \rangle$. Letting $v_+(t) = \frac{1}{\sqrt{2}}(w_\beta(t) + a(t))$ and $v_-(t) = \frac{1}{\sqrt{2}}(w_\beta(t) - a(t))$, we find the same decoupled stochastic processes as in Appendix I.3.

$$\frac{d}{dt} v_+(t) = \gamma (\beta_\gamma - \beta(t)) v_+(t), \quad \frac{d}{dt} v_-(t) = -\gamma (\beta_\gamma - \beta(t)) v_-(t) \quad (130)$$

We can use these equations to perform the necessary averages for κ_γ and D_γ . Lastly, we use

$$\frac{\partial}{\partial \beta_\gamma(s)} w_\gamma(t) = -\gamma \Theta(t-s) a(s) \quad (131)$$

to evaluate $D_\gamma(t, s)$. The observed covariances are just

$$\boldsymbol{\Sigma}_\gamma = (\gamma \mathbf{I} - \mathbf{D}_\gamma)^{-1}, \quad \boldsymbol{\Sigma}_\beta = (\gamma \mathbf{I} - \mathbf{D}_\beta)^{-1} \quad (132)$$

We note that these expressions are identical to those in Appendix J under the substitution $\beta_\gamma - \beta(t) \rightarrow \Delta(t)$ and $D \rightarrow P$. Thus the expected test risk is

$$| \langle \mathbf{w}(t) - \boldsymbol{\gamma} \rangle|^2 \sim (\beta(t) - \beta_\gamma)^2 + \frac{1}{N} \boldsymbol{\Sigma}_\gamma(t, t) + \frac{(D-1)}{N} \boldsymbol{\Sigma}_\beta(t, t) + \mathcal{O}(N^{-2}) \quad (133)$$

This recovers the variance we obtained in the multiple-sample whitened data case J.

We can then automatically differentiate the DMFT action to get the propagator. For example, for a three layer linear network, the full DMFT action has the form

$$\begin{aligned}
S = & \frac{1}{2} \text{Tr} \begin{pmatrix} \hat{H}^1 & 0 & \mathbf{I} & -\mathbf{D}^1 \\ 0 & -\hat{H}^1 & -\hat{C}^1 & \mathbf{I} \\ \mathbf{I} & -\hat{C}^1 & \mathbf{1}\mathbf{1}^T & 0 \\ -\mathbf{D}^1 & \mathbf{I} & 0 & \mathbf{G}^2 \end{pmatrix} - \gamma^2 \text{Tr} \mathbf{A} \mathbf{B} \\
& - \frac{1}{2} \ln \det \begin{pmatrix} \hat{H}^1 & 0 & \mathbf{I} & -\mathbf{D}^1 \\ 0 & -\hat{H}^1 & -\hat{C}^1 & \mathbf{I} \\ \mathbf{I} & -\hat{C}^1 & \mathbf{1}\mathbf{1}^T & 0 \\ -\mathbf{D}^1 & \mathbf{I} & 0 & \mathbf{G}^2 \end{pmatrix} \\
& - \frac{1}{2} \ln \det \begin{pmatrix} \hat{H}^2 & 0 & \mathbf{I} & -\mathbf{D}^2 \\ 0 & -\hat{H}^2 & -\hat{C}^2 & \mathbf{I} \\ \mathbf{I} & -\hat{C}^2 & \mathbf{H}^1 & 0 \\ -\mathbf{D}^2 & \mathbf{I} & 0 & \mathbf{1}\mathbf{1}^T \end{pmatrix}
\end{aligned} \tag{136}$$

where $\mathbf{C}^1 = \gamma \odot \mathbf{H}^1 + \gamma \mathbf{A}$ and $\mathbf{D}^1 = \gamma \odot \mathbf{G}^2 + \gamma \mathbf{B}$ and $\mathbf{D}^2 = \gamma \odot \mathbf{H}^2 + \gamma \mathbf{A}$. This above example can be extended to deeper networks. The total size of the block matrices which we compute determinants over is $4PT \times 4PT$ for a dataset of size P trained for T steps.

M Discrete Time Dynamics and Edge of Stability Effects

Large step size effects can induce qualitatively different dynamics in neural network training. For instance, if the step size exceeds that required for linear stability with the initial kernel, the kernel can decrease in order to stabilize the dynamics [57]. Alternatively, during training the kernel may exhibit a ‘‘progressive sharpening’’ phase where its top eigenvalue grows before reaching a stability bound set by the learning rate [19]. It is therefore well motivated to study how dynamics in this regime alter finite size effects in neural networks. We will first solve a special model which was considered in prior work [57]: a two layer linear network trained on a single training point. We will then provide the full DMFT equations for the discrete time case and provide an outline for how one could obtain finite size effects in that picture.

M.1 Two Layer Linear Equations

In a two layer linear network, the DMFT equations are

$$\begin{aligned}
h(t+1) &= h(t) + \eta \gamma \Delta(t) z(t), \quad z(t+1) = z(t) + \eta \gamma \Delta(t) h(t) \\
f(t) &= \frac{1}{\gamma} \langle z(t) h(t) \rangle
\end{aligned} \tag{137}$$

The NTK has the form $K(t) = h(t)^2 + z(t)^2$. We can easily show that the kernel and error have coupled dynamics

$$\begin{aligned}
f(t+1) &= f(t) + \eta \langle h(t)^2 + z(t)^2 \rangle \Delta(t) + \eta^2 \gamma \Delta(t)^2 \langle h(t) z(t) \rangle \\
&= f(t) + \eta K(t) \Delta(t) + \eta^2 \gamma^2 \Delta(t)^2 f(t)
\end{aligned} \tag{138}$$

$$\begin{aligned}
K(t+1) &= K(t) + 4\eta \gamma \Delta(t) \langle h(t) z(t) \rangle + \eta^2 \gamma^2 \Delta(t)^2 \langle h(t)^2 + z(t)^2 \rangle \\
&= K(t) + 4\eta \gamma^2 \Delta(t) f(t) + \eta^2 \gamma^2 \Delta(t)^2 K(t)
\end{aligned} \tag{139}$$

These equations define the infinite width evolution of $\Delta(t)$ and $K(t)$. Already at this level of analysis, we can reason about the evolution of $K(t)$. In the small η limit, we could disregard terms of order $\mathcal{O}(\eta^2)$ and arrive at the following gradient flow approximation for $K(t) \sim 2 \frac{1}{1 + \gamma^2 f(t)^2}$ [9]. This evolution will not reach the edge of stability provided that $\eta < \frac{1}{\sqrt{1 + 2\gamma^2}}$. For large γ and $y = 1$, this leads to the constraint $\eta \gamma < 1$. However, if η exceeds this bound, the gradient flow approximation is no longer reasonable and the system reaches an edge of stability effect as shown in Figure 6.

To calculate the finite size effects, we need to compute κ and $D(t, s) = \frac{\partial}{\partial (s)} \langle h(t)^2 + z(t)^2 \rangle$. To evaluate these quantities we utilize the same change of variables employed in Appendix I.3. In discrete time, these decoupled equations are

$$v_+(t+1) = v_+(t) + \eta \gamma \Delta(t) v_+(t), \quad v_-(t+1) = v_-(t) - \eta \gamma \Delta(t) v_-(t). \tag{140}$$

Given $\Delta(t)$, these can be expressed as linear systems of equations. Now, we can easily compute the uncoupled kernel variance

$$\begin{aligned}\kappa(t, s) &= 2 \langle h(t)h(s) \rangle^2 + 2 \langle z(t)z(s) \rangle^2 + 2 \langle h(t)z(s) \rangle^2 + 2 \langle z(t)h(s) \rangle^2 \\ &= \langle v_+(t)v_+(s) + v_-(t)v_-(s) \rangle^2 + \langle v_+(t)v_+(s) - v_-(t)v_-(s) \rangle^2.\end{aligned}\quad (141)$$

Similarly, we can calculate $D(t, s)$ by using the fact $h(t)^2 + z(t)^2 = v_+(t)^2 + v_-(t)^2$

$$\begin{aligned}D(t, s) &= 2 v_+(t) \frac{\partial v_+(t)}{\partial \Delta(s)} + 2 v_-(t) \frac{\partial v_-(t)}{\partial \Delta(s)} \\ \frac{\partial v_+(t)}{\partial \Delta(s)} &= \gamma \Theta(t-s) v_+(s) + \prod_{t^\theta < t} \Delta(t^\theta) \frac{\partial v_+(t^\theta)}{\partial \Delta(s)} \\ \frac{\partial v_-(t)}{\partial \Delta(s)} &= -\gamma \Theta(t-s) v_-(s) - \prod_{t^\theta < t} \Delta(t^\theta) \frac{\partial v_-(t^\theta)}{\partial \Delta(s)}\end{aligned}\quad (142)$$

These can be directly solved as a linear system of equations.

N Computing Details

Experiments for Figures 3, 6 and 2 were conducted on a Google Colab GPU with JAX. Experiments for Figures 5, A.3, 7 were performed on a NVIDIA SMX4-A100-80GB GPU. The total compute required for all Figures in the paper took around 4 hours. Jupyter Notebooks to reproduce plots can be found at https://github.com/Pehlevan-Group/dmft_fluctuations.