

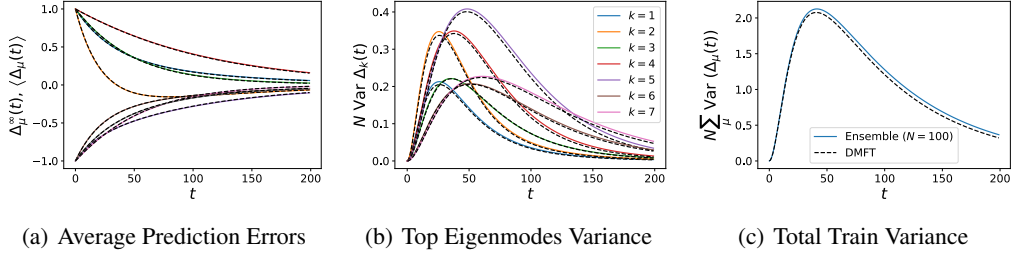
472 **A Additional Figures**

Figure A.1: We show the accuracy of the lazy-limit ODE in equation (where) 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  $\Delta_k(t) = \Delta(t) \cdot \phi_k$ . These are not ordered simply by magnitude of eigenvalues or the target projections  $y_k = \mathbf{y} \cdot \phi_k$ , but rather depend on all eigenvalue gaps  $\lambda_k - \lambda_\ell$  for  $k \neq \ell$  and also the  $\kappa_{k\ell nm}$  tensor. (c) The total variance for all training points  $N \sum_\mu \text{Var} \Delta_\mu(t) = N \sum_k \text{Var} \Delta_k(t)$  is also well predicted by the DMFT propagator equations.

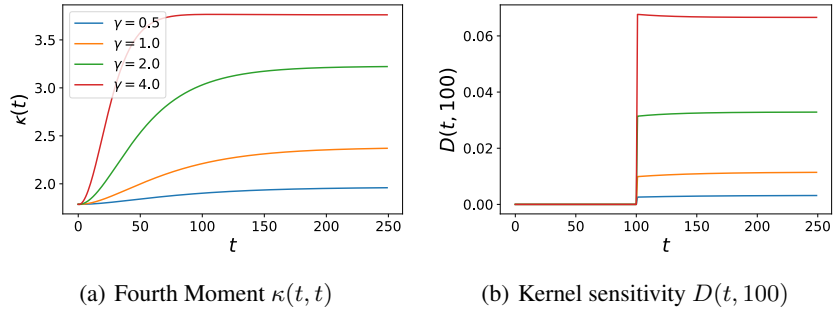


Figure A.2: The  $\kappa$  and  $D$  functions for varying  $\gamma$  in Figure 2. (a) The uncoupled kernel variance  $\kappa(t, t)$  increases monotonically with  $\gamma$ . This reveals that the dynamical filtering of  $\kappa$  is what is responsible for the late time decrease in variance during feature learning. (b) The tensor  $D(t, s)$  measures sensitivity of kernel at time  $t$  to perturbation in  $\Delta$  at time  $s$ . The  $D(t, s)$  entries also increase with  $\gamma$ . This suggests that the reduction in variance of the training error and the kernel are not due to reduction in  $\kappa$ , but rather a dynamical filtering effect due to scale growth in  $K_\infty$  and rapid reduction in  $\Delta_\infty$ .

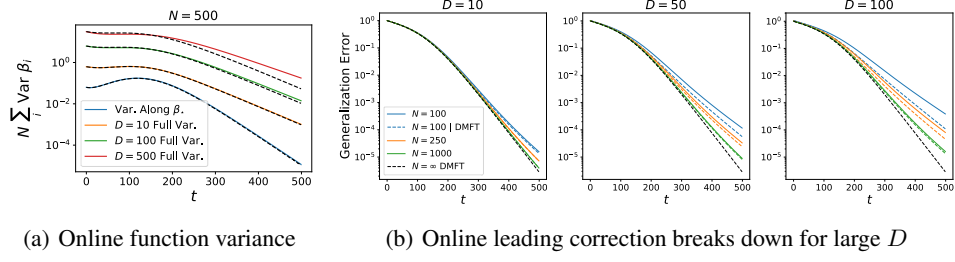


Figure A.3: Online learning follows identical finite size effects to offline training in two layer linear networks. (a) The variance of  $\beta(t)$  in online learning vs input dimension  $D$ . (b) The predicted leading correction to the generalization error is accurate for  $N \gg D$  but breaks down as  $D$  becomes comparable to  $N$ .

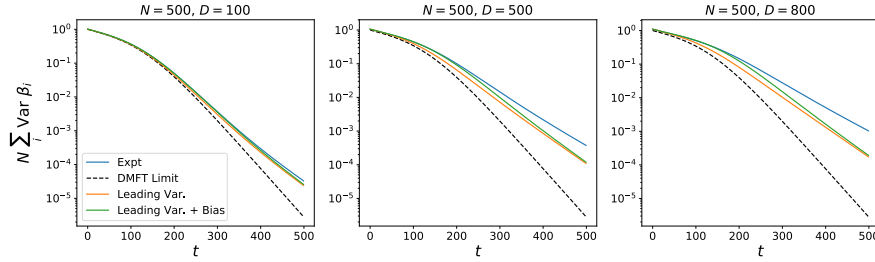


Figure A.4: A comparison of the bias and variance corrections in the toy model of Figure 3. At small  $D/N$  (or  $P/N$  for offline training) the leading variance and the leading variance and leading bias both track the experiment. Both the bias and the variance contribute positively towards the total 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  $D^2/N^2$ ,  $D^3/N^3$ , ...).

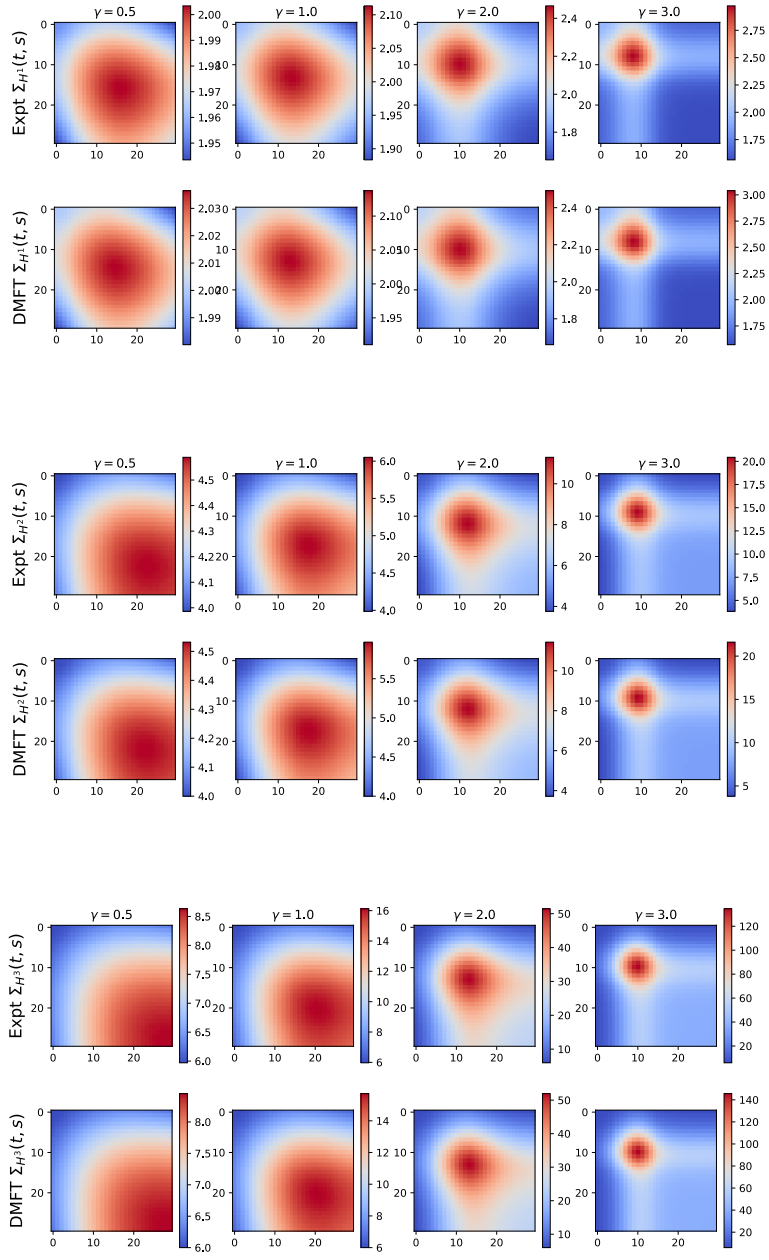


Figure A.5: The covariance of kernel entries across pairs of time points  $\Sigma_{H^\ell}(t, s) = N \text{Cov}(H^\ell(t, t), H^\ell(s, s))$  for depth 4 linear network trained on whitened data. The variance becomes increasingly localized in time as feature learning  $\gamma$  increases.

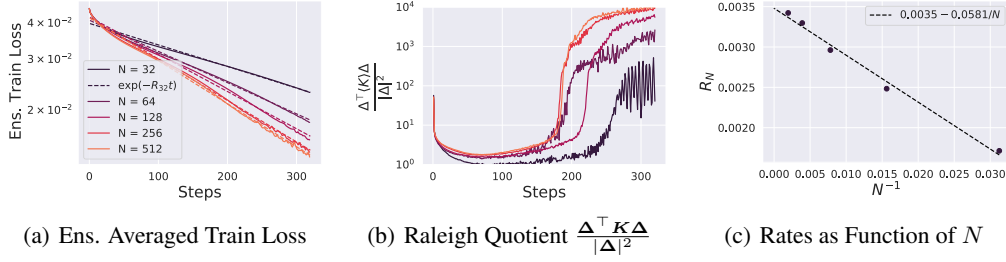


Figure A.6: The ensemble averaged train loss for the same depth 6 CNN trained on a random subsample of  $P = 64$  CIFAR-10 points. Training is full batch gradient descent with  $\gamma = 0.05$ . (a) The ensemble train accuracy for this subset of CIFAR-10 is well modeled as an exponential in time  $\mathcal{L}(t)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

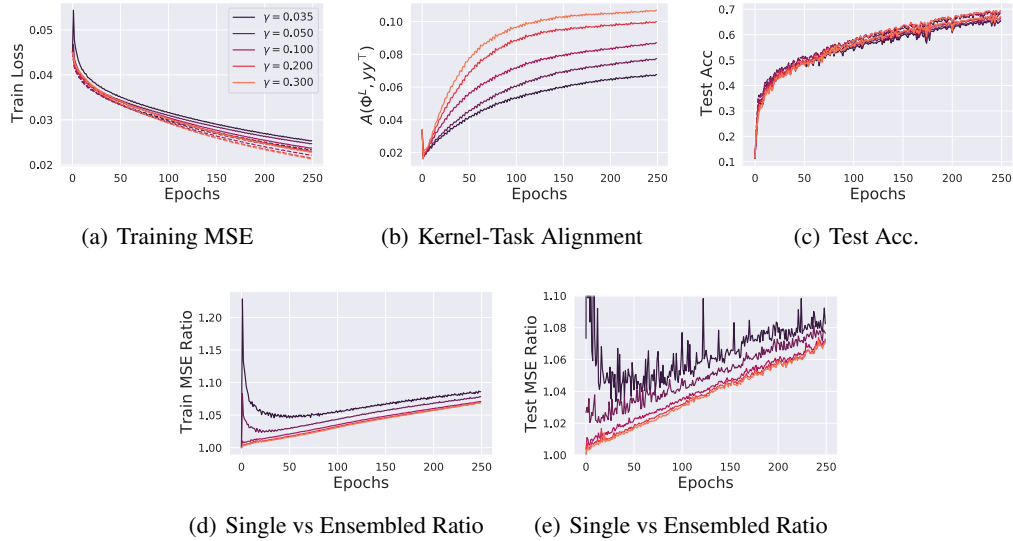


Figure A.7: Width  $N = 64$  depth 6 CNNs trained on the full CIFAR-10 with MSE. An ensemble of size  $E = 10$  randomly initialized networks are trained. (a) Training MSE for varying  $\gamma$ . (b) Final layer kernel-task alignment does strongly depend on  $\gamma$ , despite similar train dynamics. (c) Top-1 classification test accuracy is only slightly different across  $\gamma$ . A small benefit from ensembleing 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  $\gamma$  as evidenced by the alignment (cosine similarity) between the final layer feature kernel  $\Phi^L$  and the target test labels  $y$ .

## 473 B CIFAR-10 Experimental Details

474 We trained the following depth 6 CNN architecture in the mean field parameterization using FLAX  
 475 [58] on a single GPU. The bias parameters were zero in each hidden Conv layer, but were used for  
 476 the readout weights. The networks were trained with MSE loss on centered 10 dimensional targets  
 477  $\mathbf{y}_\mu \in \mathbb{R}^{10}$  for  $\mu \in [P]$ . Each convolution was followed by an average pooling operation. To obtain  
 478 mean field behavior, NTK parameterization with a modified final layer is used [7, 9].

```

479 1 from flax import linen as nn
480 2 import jax.numpy as jnp
481 3
482 4 class CNN(nn.Module):
483 5
484 6     width: int
485 7
486 8     def setup(self):
487 9         kif = nn.initializers.normal(stddev = 1.0) # 0_N(1) entries
488 10        self.conv1 = nn.Conv(features = self.width, kernel_init = kif,
489        use_bias = False, kernel_size = (3,3))
490 11        self.conv2 = nn.Conv(features = self.width, kernel_init = kif,
491        use_bias = False, kernel_size = (3,3))
492 12        self.conv3 = nn.Conv(features = self.width, kernel_init = kif,
493        use_bias = False, kernel_size = (3,3))
494 13        self.conv4 = nn.Conv(features = self.width, kernel_init = kif,
495        use_bias = False, kernel_size = (3,3))
496 14        self.conv5 = nn.Conv(features = self.width, kernel_init = kif,
497        use_bias = False, kernel_size = (3,3))
498 15        self.readout = nn.Dense(features = 10, use_bias = True,
499        kernel_init = kif)
500 16        return
501 17
502 18    def __call__(self, x, train = True):
503 19        N = self.width
504 20        D = 3
505 21        x = self.conv1(x) / jnp.sqrt(D * 9)
506 22        x = jnp.sqrt(2.0) * nn.relu(x)
507 23        x = nn.avg_pool(x, window_shape=(2,2), strides = (2,2)) # 32 x 32
508        -> 16 x 16
509 24        x = self.conv2(x) / jnp.sqrt(N*9) # explicit N^{-1/2}
510 25        x = jnp.sqrt(2.0) * nn.relu(x)
511 26        x = nn.avg_pool(x, window_shape=(2,2), strides = (2,2)) # 16 x 16
512        -> 8 x 8
513 27        x = self.conv3(x) / jnp.sqrt(N*9)
514 28        x = jnp.sqrt(2.0) * nn.relu(x)
515 29        x = nn.avg_pool(x, window_shape=(2,2), strides = (2,2)) # 8 x 8 ->
516        4 x 4
517 30        x = self.conv4(x) / jnp.sqrt(N*9)
518 31        x = jnp.sqrt(2.0) * nn.relu(x)
519 32        x = nn.avg_pool(x, window_shape=(2,2), strides = (2,2)) # 4 x 4
520        -> 2 x 2
521 33        x = self.conv5(x) / jnp.sqrt(N*9)
522 34        x = jnp.sqrt(2.0) * nn.relu(x)
523 35        x = nn.avg_pool(x, window_shape=(2,2), strides = (2,2)) # 2 x 2
524        -> 1 x 1
525 36        x = x.reshape((x.shape[0], -1)) # flatten
526 37        x = self.readout(x) / N # for mean field scaling
527 38        return x

```

528 All models were trained with standard SGD with a batch size of 256. Each element in the ensemble of  
 529  $E$  networks is trained on identical batches presented in identical order. For the Figure 6 experiments,  
 530 the raw learning rate is scaled as  $\eta = 10N\sqrt{\gamma}$  with  $\gamma = 0.2$  (note that mean field theory requires  
 531 scaling the raw learning rate linearly with  $N$  since the raw NTK is  $\mathcal{O}(N^{-1})$  [9]). For Figure A.7, the  
 532 learning rate is  $\eta = 5N\sqrt{\gamma}$ . We find that choosing  $\eta \propto \sqrt{\gamma}$  gives approximately conserved training

times across  $\gamma$  (though distinct representation dynamics). The Figure [A.6](#) shows the dynamics of fitting  $P = 64$  training points with full batch gradient descent and  $\gamma = 0.1$ .

## C 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

$$\langle O(\mathbf{q}) \rangle_N = \frac{\int d\mathbf{q} \exp(N S(\mathbf{q})) O(\mathbf{q})}{\int d\mathbf{q} \exp(N S(\mathbf{q}))} \quad (11)$$

As discussed in the main text, the  $N \rightarrow \infty$  limit gives  $\langle O(\mathbf{q}) \rangle_N \sim O(\mathbf{q}_\infty)$  where  $\frac{\partial S}{\partial \mathbf{q}}|_{\mathbf{q}_\infty} = 0$  by a steepest descent argument [\[55\]](#). We assume that  $S$ 's Hessian is negative semidefinite so that  $\Sigma \equiv -[\nabla^2 S(\mathbf{q})|_{\mathbf{q}_\infty}]^{-1} \succeq 0$  and Taylor expand  $S(\mathbf{q})$  around the saddle point  $\mathbf{q}_\infty$  giving  $S(\mathbf{q}) = S(\mathbf{q}_\infty) + \frac{1}{2}(\mathbf{q} - \mathbf{q}_\infty)^\top \nabla^2 S(\mathbf{q})(\mathbf{q} - \mathbf{q}_\infty) + V(\mathbf{q} - \mathbf{q}_\infty)$ . We note that the remainder function  $V$  contains only cubic and higher powers of  $\mathbf{q} - \mathbf{q}_\infty \equiv \delta/\sqrt{N}$ . The variable  $\delta$  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

$$\begin{aligned} \langle O(\mathbf{q}) \rangle_N &= \frac{\int d\mathbf{q} \exp\left(-\frac{N}{2}(\mathbf{q} - \mathbf{q}_\infty)^\top \Sigma^{-1}(\mathbf{q} - \mathbf{q}_\infty) + NV(\mathbf{q} - \mathbf{q}_\infty)\right) O(\mathbf{q})}{\int d\mathbf{q} \exp\left(-\frac{N}{2}(\mathbf{q} - \mathbf{q}_\infty)^\top \Sigma^{-1}(\mathbf{q} - \mathbf{q}_\infty) + NV(\mathbf{q} - \mathbf{q}_\infty)\right)} \\ &= \frac{\int d\delta \exp\left(-\frac{1}{2}\delta^\top \Sigma^{-1}\delta\right) (1 + NV + \frac{N^2}{2}V^2 + \dots) O(\mathbf{q}_\infty + N^{-1/2}\delta)}{\int d\delta \exp\left(-\frac{1}{2}\delta^\top \Sigma^{-1}\delta\right) (1 + NV + \frac{N^2}{2}V^2 + \dots)} \\ &= \frac{\langle O \rangle_\infty + N \langle VO \rangle_\infty + \frac{N^2}{2!} \langle V^2 O \rangle_\infty + \frac{N^3}{3!} \langle V^3 O \rangle_\infty + \dots}{1 + N \langle V \rangle_\infty + \frac{N^2}{2!} \langle V^2 \rangle_\infty + \frac{N^3}{3!} \langle V^3 \rangle_\infty + \dots} \\ &= \langle O \rangle_\infty \frac{1 + N \langle VO \rangle_\infty / \langle O \rangle_\infty + \frac{N^2}{2!} \langle V^2 O \rangle_\infty / \langle O \rangle_\infty + \frac{N^3}{3!} \langle V^3 O \rangle_\infty / \langle O \rangle_\infty + \dots}{1 + N \langle V \rangle_\infty + \frac{N^2}{2!} \langle V^2 \rangle_\infty + \frac{N^3}{3!} \langle V^3 \rangle_\infty + \dots} \end{aligned} \quad (12)$$

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

$$\frac{N^k \langle V^k O \rangle_\infty}{\langle O \rangle_\infty}, N^k \langle V^k \rangle_\infty \sim \begin{cases} \mathcal{O}(N^{-(k+1)/2}) & k \text{ odd} \\ \mathcal{O}(N^{-k/2}) & k \text{ even} \end{cases} \quad (13)$$

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  $\frac{\langle O \rangle}{\langle O \rangle_\infty}$  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

$$\begin{aligned} \langle O \rangle_N &= \langle O \rangle_\infty + N[\langle OV \rangle_\infty - \langle O \rangle_\infty \langle V \rangle_\infty] \\ &\quad + \frac{N^2}{2} [\langle V^2 O \rangle_\infty - 2 \langle VO \rangle_\infty \langle V \rangle_\infty + 2 \langle V \rangle_\infty^2 \langle O \rangle_\infty - \langle V^2 \rangle_\infty \langle O \rangle_\infty] + \dots \end{aligned} \quad (14)$$

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  $O$ .

### 557 C.1 Square Deviation from DMFT

558 We will now analyze the fluctuation statistics of our order parameters around the saddle point  
 559  $\langle (\mathbf{q} - \mathbf{q}_\infty)(\mathbf{q} - \mathbf{q}_\infty)^\top \rangle_N$  which has the form

$$\begin{aligned} \langle (\mathbf{q} - \mathbf{q}_\infty)(\mathbf{q} - \mathbf{q}_\infty)^\top \rangle_N &= \frac{\langle (\mathbf{q} - \mathbf{q}_\infty)(\mathbf{q} - \mathbf{q}_\infty)^\top \rangle_\infty + N \langle V (\mathbf{q} - \mathbf{q}_\infty)(\mathbf{q} - \mathbf{q}_\infty)^\top \rangle_\infty + \dots}{1 + N \langle V \rangle_\infty + \dots} \\ &= \left[ \frac{\frac{1}{N} \Sigma + \mathcal{O}(N^{-2})}{1 + \mathcal{O}(N^{-1})} \right] \sim \frac{1}{N} \Sigma + \mathcal{O}(N^{-2}), \end{aligned} \quad (15)$$

560 as stated in the main text and verified empirically in Figure 2(a). The reason that the terms in the  
 561 numerator involving  $V$  can be no larger than  $\mathcal{O}(N^{-2})$  comes from vanishing of odd moments for  
 562  $\mathbf{q} - \mathbf{q}_\infty$  in the unperturbed distribution. Thus the leading expression for  $\langle (\mathbf{q} - \mathbf{q}_\infty)(\mathbf{q} - \mathbf{q}_\infty)^\top \rangle$  only  
 563 depends on  $\Sigma$  and not on  $V$ .

### 564 C.2 Mean Deviation from DMFT

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

$$\begin{aligned} \langle \mathbf{q} - \mathbf{q}_\infty \rangle_N &= \frac{\langle \mathbf{q} - \mathbf{q}_\infty \rangle_\infty + N \langle (\mathbf{q} - \mathbf{q}_\infty) V \rangle_\infty + \frac{N^2}{2} \langle (\mathbf{q} - \mathbf{q}_\infty) V^2 \rangle_\infty + \dots}{1 + N \langle V \rangle_\infty + \frac{N^2}{2} \langle V^2 \rangle_\infty + \dots} \\ &\sim \frac{\Sigma \left\langle \frac{\partial V}{\partial \mathbf{q}} \right\rangle_\infty + \mathcal{O}(N^{-2})}{1 + \mathcal{O}(N^{-1})} \sim \Sigma \left\langle \frac{\partial V}{\partial \mathbf{q}} \right\rangle_\infty + \mathcal{O}(N^{-2}). \end{aligned} \quad (16)$$

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

### 573 C.3 Covariance of Order Parameters

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

$$\begin{aligned} \text{Cov}(\mathbf{q}) &= \left\langle (\mathbf{q} - \langle \mathbf{q} \rangle_N) (\mathbf{q} - \langle \mathbf{q} \rangle_N)^\top \right\rangle_N \\ &= \left\langle (\mathbf{q} - \mathbf{q}_\infty) (\mathbf{q} - \mathbf{q}_\infty)^\top \right\rangle_N - \left\langle (\mathbf{q}_\infty - \langle \mathbf{q} \rangle_N) (\mathbf{q}_\infty - \langle \mathbf{q} \rangle_N)^\top \right\rangle_N \\ &\sim \frac{1}{N} \Sigma + \mathcal{O}(N^{-2}) \end{aligned} \quad (17)$$

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

## 580 D Propagator Structure for the full DMFT Action

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

587 pre-gradient fields. In this section, we will mainly focus on gradient flow, but we describe large step  
 588 size in Appendix [E](#).

$$\begin{aligned}
 S = & \sum_{\ell\mu\nu} \int dt ds \left[ \hat{\Phi}_{\mu\nu}^{\ell}(t, s) \Phi_{\mu\nu}^{\ell}(t, s) + \hat{G}_{\mu\nu}^{\ell}(t, s) G_{\mu\nu}^{\ell}(t, s) - \gamma^2 A_{\nu\mu}^{\ell}(s, t) B_{\mu\nu}^{\ell}(t, s) \right] \\
 & + \sum_{\mu} \int dt \hat{\Delta}_{\mu}(t) \left[ \Delta_{\mu}(t) - y_{\mu} + \sum_{\nu} \int ds \Theta(t-s) K_{\mu\nu}(s) \Delta_{\nu}(s) \right] \\
 & + \sum_{\mu\nu} \int dt \hat{K}_{\mu\nu}(t) \left[ K_{\mu\nu}(t) - \sum_{\ell} G_{\mu\nu}^{\ell+1}(t) \Phi_{\mu\nu}^{\ell}(t) \right] \\
 & + \sum_{\ell} \ln \mathcal{Z}_{\ell}[\Delta, \hat{\Phi}^{\ell}, \hat{G}^{\ell}, \Phi^{\ell-1}, G^{\ell+1}, A^{\ell-1}, B^{\ell}]
 \end{aligned} \tag{18}$$

589 where the single site moment generating functionals  $\mathcal{Z}_{\ell}$  have the form

$$\begin{aligned}
 \mathcal{Z}_{\ell} = & \mathbb{E}_{\{h_{\mu}^{\ell}(t), z_{\mu}^{\ell}(t)\}} \exp \left( - \sum_{\mu\nu} \int dt ds \left[ \phi(h_{\mu}^{\ell}(t)) \phi(h_{\nu}^{\ell}(s)) \hat{\Phi}_{\mu\nu}^{\ell}(t, s) + g_{\mu}^{\ell}(t) g_{\nu}^{\ell}(s) \hat{G}_{\mu\nu}^{\ell}(t, s) \right] \right) \\
 h_{\mu}^{\ell}(t) = & u_{\mu}^{\ell}(t) + \gamma \int_0^t ds \sum_{\nu} [\Phi_{\mu\nu}^{\ell-1}(t, s) \Delta_{\nu}(s) + A_{\mu\nu}^{\ell-1}(t, s)] g_{\nu}^{\ell}(s), \quad \{u_{\mu}^{\ell}(t)\} \sim \mathcal{GP}(0, \Phi^{\ell-1}) \\
 z_{\mu}^{\ell}(t) = & r_{\mu}^{\ell}(t) + \gamma \int_0^t ds \sum_{\nu} [G_{\mu\nu}^{\ell+1}(t, s) \Delta_{\nu}(s) + B_{\mu\nu}^{\ell}(t, s)] \phi(h_{\nu}^{\ell}(s)), \quad \{r_{\mu}^{\ell}(t)\} \sim \mathcal{GP}(0, G^{\ell+1})
 \end{aligned} \tag{19}$$

590 with  $g_{\mu}^{\ell}(t) = \dot{\phi}(h_{\mu}^{\ell}(t)) z_{\mu}^{\ell}(t)$ . The saddle point equations give the infinite width evolution of our order  
 591 parameters.

$$\begin{aligned}
 \frac{\partial S}{\partial \hat{\Phi}_{\mu\nu}^{\ell}(t, s)} &= \Phi_{\mu\nu}^{\ell}(t, s) - \langle \phi(h_{\mu}^{\ell}(t)) \phi(h_{\nu}^{\ell}(s)) \rangle = 0 \\
 \frac{\partial S}{\partial \hat{G}_{\mu\nu}^{\ell}(t, s)} &= G_{\mu\nu}^{\ell}(t, s) - \langle g_{\mu}^{\ell}(t) g_{\nu}^{\ell}(s) \rangle = 0 \\
 \frac{\partial S}{\partial A_{\nu\mu}^{\ell}(s, t)} &= -\gamma^2 B_{\mu\nu}^{\ell}(t, s) + \gamma \left\langle \frac{\partial \phi(h_{\mu}^{\ell}(t))}{\partial r_{\nu}^{\ell}(s)} \right\rangle = 0 \\
 \frac{\partial S}{\partial B_{\nu\mu}^{\ell}(s, t)} &= -\gamma^2 A_{\mu\nu}^{\ell}(t, s) + \gamma \left\langle \frac{\partial g_{\mu}^{\ell}(t)}{\partial u_{\nu}^{\ell}(s)} \right\rangle = 0 \\
 \frac{\partial S}{\partial \hat{K}_{\mu\nu}(t)} &= K_{\mu\nu}(t) - \sum_{\ell} G_{\mu\nu}^{\ell+1}(t, t) \Phi_{\mu\nu}^{\ell}(t, t) = 0 \\
 \frac{\partial S}{\partial \hat{\Delta}_{\mu}(t)} &= \Delta_{\mu}(t) - y_{\mu} + \int_0^t ds \sum_{\nu} K_{\mu\nu}(s) \Delta_{\nu}(s) = 0
 \end{aligned} \tag{20}$$

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



597 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}_{\mu\nu}^\ell(t, s) \partial \hat{\Phi}_{\alpha\beta}^\ell(t', s')} &= \langle \phi(h_\mu^\ell(t)) \phi(h_\nu^\ell(s)) \phi(h_\alpha^\ell(t')) \phi(h_\beta^\ell(s')) \rangle - \Phi_{\mu\nu}^\ell(t, s) \Phi_{\alpha\beta}^\ell(t', s') \\
&\equiv \kappa_{\mu\nu\alpha\beta}^{\Phi^\ell}(t, s, t', s') \\
\frac{\partial^2 S}{\partial \hat{G}_{\mu\nu}^\ell(t, s) \partial \hat{G}_{\alpha\beta}^\ell(t', s')} &= \langle g_\mu^\ell(t) g_\nu^\ell(s) g_\alpha^\ell(t') g_\beta^\ell(s') \rangle - G_{\mu\nu}^\ell(t, s) G_{\alpha\beta}^\ell(t', s') \\
&\equiv \kappa_{\mu\nu\alpha\beta}^{G^\ell}(t, s, t', s') \\
\frac{\partial^2 S}{\partial \hat{\Phi}_{\mu\nu}^\ell(t, s) \partial \hat{G}_{\alpha\beta}^\ell(t', s')} &= \langle \phi(h_\mu^\ell(t)) \phi(h_\nu^\ell(s)) g_\alpha^\ell(t') g_\beta^\ell(s') \rangle - \Phi_{\mu\nu}^\ell(t, s) G_{\alpha\beta}^\ell(t', s') \\
&\equiv \kappa_{\mu\nu\alpha\beta}^{\Phi^\ell G^\ell}(t, s, t', s') \\
\frac{\partial^2 S}{\partial \hat{\Phi}_{\mu\nu}^\ell(t, s) \partial A_{\beta\alpha}^{\ell-1}(s', t')} &= -\gamma \left\langle \frac{\partial \phi(h_\mu^\ell(t))}{\partial u_\beta^\ell(s')} \phi(h_\nu^\ell(s)) g_\alpha^\ell(t') \right\rangle \\
&\quad - \gamma \left\langle \phi(h_\mu^\ell(t)) \frac{\partial \phi(h_\nu^\ell(t))}{\partial u_\beta^\ell(s')} g_\alpha^\ell(t') \right\rangle \\
&\quad - \gamma \left\langle \phi(h_\mu^\ell(t)) \phi(h_\nu^\ell(s)) \frac{\partial g_\alpha^\ell(t')}{\partial u_\beta^\ell(s')} \right\rangle - \gamma^2 \Phi_{\mu\nu}^\ell(t, s) B_{\alpha\beta}^{\ell-1}(t', s') \\
&\equiv -\gamma \kappa_{\mu\nu\alpha\beta}^{\Phi^\ell B^{\ell-1}}(t, s) \\
\frac{\partial^2 S}{\partial \hat{\Phi}_{\mu\nu}^\ell(t, s) \partial B_{\beta\alpha}^\ell(s', t')} &= -\gamma \left\langle \frac{\partial \phi(h_\mu^\ell(t))}{\partial r_\beta^\ell(s')} \phi(h_\nu^\ell(s)) \phi(h_\alpha^\ell(t')) \right\rangle \\
&\quad - \gamma \left\langle \phi(h_\mu^\ell(t)) \frac{\partial \phi(h_\nu^\ell(t))}{\partial r_\beta^\ell(s')} \phi(h_\alpha^\ell(t')) \right\rangle \\
&\quad - \gamma \left\langle \phi(h_\mu^\ell(t)) \phi(h_\nu^\ell(s)) \frac{\partial \phi(h_\alpha^\ell(t'))}{\partial r_\beta^\ell(s')} \right\rangle - \gamma^2 \Phi_{\mu\nu}^\ell(t, s) A_{\alpha\beta}^\ell(t', s') \\
&\equiv -\gamma \kappa_{\mu\nu\alpha\beta}^{\Phi^\ell A^\ell}(t, s) \\
\frac{\partial^2 S}{\partial \hat{G}_{\mu\nu}^\ell(t, s) \partial A_{\beta\alpha}^{\ell-1}(s', t')} &= -\gamma \left\langle \frac{\partial g_\mu^\ell(t)}{\partial u_\beta^\ell(s')} g_\nu^\ell(s) g_\alpha^\ell(t') \right\rangle - \gamma \left\langle g_\mu^\ell(t) \frac{\partial g_\nu^\ell(t)}{\partial u_\beta^\ell(s')} g_\alpha^\ell(t') \right\rangle \\
&\quad - \gamma \left\langle g_\mu^\ell(t) g_\nu^\ell(s) \frac{\partial g_\alpha^\ell(t')}{\partial u_\beta^\ell(s')} \right\rangle - \gamma^2 G_{\mu\nu}^\ell(t, s) B_{\alpha\beta}^{\ell-1}(t', s') \\
&\equiv -\gamma \kappa_{\mu\nu\alpha\beta}^{G^\ell B^{\ell-1}}(t, s) \\
\frac{\partial^2 S}{\partial \hat{G}_{\mu\nu}^\ell(t, s) \partial B_{\beta\alpha}^\ell(s', t')} &= -\gamma \left\langle \frac{\partial g_\mu^\ell(t)}{\partial r_\beta^\ell(s')} g_\nu^\ell(s) \phi(h_\alpha^\ell(t')) \right\rangle - \gamma \left\langle g_\mu^\ell(t) \frac{\partial g_\nu^\ell(t)}{\partial r_\beta^\ell(s')} \phi(h_\alpha^\ell(t')) \right\rangle \\
&= -\gamma \left\langle g_\mu^\ell(t) g_\nu^\ell(s) \frac{\partial \phi(h_\alpha^\ell(t'))}{\partial r_\beta^\ell(s')} \right\rangle - \gamma^2 G_{\mu\nu}^\ell(t, s) A_{\alpha\beta}^\ell(t', s') \\
&\equiv -\gamma \kappa_{\mu\nu\alpha\beta}^{G^\ell A^\ell}(t, s) \\
\frac{\partial^2 S}{\partial A_{\mu\nu}^\ell(t, s) \partial B_{\beta\alpha}^\ell(s', t')} &= -\gamma^2 \delta_{\mu\alpha} \delta_{\nu\beta} \delta(t - t') \delta(s - s') \\
\frac{\partial^2 S}{\partial A_{\nu\mu}^{\ell-1}(s, t) \partial B_{\beta\alpha}^\ell(s', t')} &= \gamma^2 \left\langle \frac{\partial^2}{\partial u_\nu^\ell(s) \partial r_\beta^\ell(s')} [g_\mu^\ell(t) \phi(h_\alpha^\ell(t'))] \right\rangle - \gamma^4 B_{\mu\nu}^{\ell-1}(t, s) A_{\alpha\beta}^\ell(t', s') \\
&\equiv \kappa_{\mu\nu\alpha\beta}^{B^{\ell-1} A^\ell}(t, s, t', s')
\end{aligned} \tag{21}$$

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

$$\begin{aligned}
 \frac{\partial^2 S}{\partial \Phi_{\mu\nu}^\ell(t, s) \partial \Phi_{\alpha\beta}^{\ell'}(t', s')} &= 0 \\
 \frac{\partial^2 S}{\partial G_{\mu\nu}^\ell(t, s) \partial G_{\alpha\beta}^{\ell'}(t', s')} &= 0 \\
 \frac{\partial^2 S}{\partial \Phi_{\mu\nu}^\ell(t, s) \partial G_{\alpha\beta}^{\ell'}(t', s')} &= 0 \\
 \frac{\partial^2 S}{\partial \Phi_{\mu\nu}^\ell(t, s) \partial K_{\alpha\beta}(s')} &= 0 \\
 \frac{\partial^2 S}{\partial G_{\mu\nu}^\ell(t, s) \partial K_{\alpha\beta}(s')} &= 0 \\
 \frac{\partial^2 S}{\partial \Phi_{\mu\nu}^\ell(t, s) \partial \Delta_\alpha(s')} &= 0 \\
 \frac{\partial^2 S}{\partial G_{\mu\nu}^\ell(t, s) \partial \Delta_\alpha(s')} &= 0 \\
 \frac{\partial^2 S}{\partial K_{\mu\nu}(t) \partial K_{\alpha\beta}(s)} &= 0 \\
 \frac{\partial^2 S}{\partial K_{\mu\nu}(t) \partial \Delta_\alpha(s)} &= 0 \\
 \frac{\partial^2 S}{\partial \Delta_\mu(t) \partial \Delta_\alpha(s)} &= 0
 \end{aligned} \tag{22}$$

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

$$\begin{aligned}
 \frac{\partial^2 S}{\partial \hat{\Phi}_{\mu\nu}^\ell(t, s) \partial \Delta_\alpha(t')} &= - \left\langle \frac{\partial}{\partial \Delta_\alpha(t')} [\phi(h_\mu^\ell(t)) \phi(h_\nu^\ell(s))] \right\rangle \equiv -D_{\mu\nu\alpha}^{\Phi^\ell \Delta}(t, s, t') \\
 \frac{\partial^2 S}{\partial \hat{G}_{\mu\nu}^\ell(t, s) \partial \Delta_\alpha(t')} &= - \left\langle \frac{\partial}{\partial \Delta_\alpha(t')} [g_\mu^\ell(t) g_\nu^\ell(s)] \right\rangle \equiv -D_{\mu\nu\alpha}^{G^\ell \Delta}(t, s, t') \\
 \frac{\partial^2 S}{\partial A_{\nu\mu}^{\ell-1}(s, t) \partial \Delta_\alpha(t')} &= \gamma \left\langle \frac{\partial}{\partial \Delta_\alpha(t') \partial u_\nu^\ell(s)} g_\mu^\ell(t) \right\rangle \equiv \gamma D_{\mu\nu\alpha}^{B^{\ell-1}, \Delta}(t, s, t') \\
 \frac{\partial^2 S}{\partial B_{\nu\mu}^\ell(s, t) \partial \Delta_\alpha(t')} &= \gamma \left\langle \frac{\partial}{\partial \Delta_\alpha(t') \partial r_\nu^\ell(s)} \phi(h_\mu^\ell(t)) \right\rangle \equiv \gamma D_{\mu\nu\alpha}^{A^\ell \Delta}(t, s, t')
 \end{aligned}$$

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

$$\begin{aligned}
\frac{\partial^2 S}{\partial \hat{\Phi}_{\mu\nu}^\ell(t, s) \partial \Phi_{\alpha\beta}^{\ell'}(t', s')} &= \delta_{\ell, \ell'} \delta_{\mu\nu} \delta(t - t') \delta(s - s') \\
&\quad - \delta_{\ell-1, \ell'} \frac{\partial}{\partial \Phi_{\alpha\beta}^{\ell-1}(t', s')} \langle \phi(h_\mu^\ell(t)) \phi(h_\nu^\ell(s)) \rangle \\
&\equiv \delta_{\ell, \ell'} \delta_{\mu\nu} \delta(t - t') \delta(s - s') - \delta_{\ell-1, \ell'} D_{\mu\nu\alpha\beta}^{\Phi^\ell \Phi^{\ell-1}}(t, s, t', s') \\
\frac{\partial^2 S}{\partial \hat{G}_{\mu\nu}^\ell(t, s) \partial G_{\alpha\beta}^{\ell'}(t', s')} &= \delta_{\ell, \ell'} \delta_{\mu\nu} \delta(t - t') \delta(s - s') - \delta_{\ell+1, \ell'} \frac{\partial}{\partial G_{\alpha\beta}^{\ell+1}(t', s')} \langle g_\mu^\ell(t) g_\nu^\ell(s) \rangle \\
&\equiv \delta_{\ell, \ell'} \delta_{\mu\nu} \delta(t - t') \delta(s - s') - \delta_{\ell-1, \ell'} D_{\mu\nu\alpha\beta}^{G^\ell G^{\ell+1}}(t, s, t', s') \\
\frac{\partial^2 S}{\partial \hat{\Phi}_{\mu\nu}^\ell(t, s) \partial G_{\alpha\beta}^{\ell+1}(t', s')} &= - \frac{\partial}{\partial G_{\alpha\beta}^{\ell+1}(t', s')} \langle \phi(h_\mu^\ell(t)) \phi(h_\nu^\ell(s)) \rangle \equiv -D_{\mu\nu\alpha\beta}^{\Phi^\ell, G^{\ell+1}}(t, s, t', s') \\
\frac{\partial^2 S}{\partial \hat{G}_{\mu\nu}^\ell(t, s) \partial \Phi_{\alpha\beta}^{\ell-1}(t', s')} &= - \frac{\partial}{\partial \Phi_{\alpha\beta}^{\ell-1}(t', s')} \langle g_\mu^\ell(t) g_\nu^\ell(s) \rangle \equiv -D_{\mu\nu\alpha\beta}^{G^\ell, \Phi^{\ell-1}}(t, s, t', s') \\
\frac{\partial^2 S}{\partial A_{\nu\mu}^{\ell-1}(s, t) \partial \Phi_{\alpha\beta}^{\ell-1}(t', s')} &= \gamma \frac{\partial}{\partial \Phi_{\alpha\beta}^{\ell-1}(t', s')} \left\langle \frac{\partial g_\mu^\ell(t)}{\partial r_\nu^\ell(s)} \right\rangle \equiv \gamma D_{\mu\nu\alpha\beta}^{B^{\ell-1}, \Phi^{\ell-1}}(t, s, t', s') \\
\frac{\partial^2 S}{\partial B_{\nu\mu}^\ell(s, t) \partial \Phi_{\alpha\beta}^{\ell-1}(t', s')} &= \gamma \frac{\partial}{\partial \Phi_{\alpha\beta}^{\ell-1}(t', s')} \left\langle \frac{\partial \phi(h_\mu^\ell(t))}{\partial u_\nu^\ell(s)} \right\rangle \equiv \gamma D_{\mu\nu\alpha\beta}^{A^\ell, \Phi^{\ell-1}}(t, s, t', s') \\
\frac{\partial^2 S}{\partial A_{\nu\mu}^{\ell-1}(s, t) \partial G_{\alpha\beta}^{\ell+1}(t', s')} &= \gamma \frac{\partial}{\partial G_{\alpha\beta}^{\ell+1}(t', s')} \left\langle \frac{\partial g_\mu^\ell(t)}{\partial r_\nu^\ell(s)} \right\rangle \equiv \gamma D_{\mu\nu\alpha\beta}^{B^{\ell-1}, G^{\ell+1}}(t, s, t', s') \\
\frac{\partial^2 S}{\partial B_{\nu\mu}^\ell(s, t) \partial G_{\alpha\beta}^{\ell+1}(t', s')} &= \gamma \frac{\partial}{\partial G_{\alpha\beta}^{\ell+1}(t', s')} \left\langle \frac{\partial \phi(h_\mu^\ell(t))}{\partial u_\nu^\ell(s)} \right\rangle \equiv \gamma D_{\mu\nu\alpha\beta}^{A^\ell, G^{\ell+1}}(t, s, t', s') \quad (23)
\end{aligned}$$

We note that terms such as  $\frac{\partial}{\partial \Phi_{\alpha\beta}^{\ell-1}(t', s')} \langle \phi(h_\mu^\ell(t)) \phi(h_\nu^\ell(s)) \rangle$  can be further decomposed since the average over the  $\{u_\mu^\ell(t)\} \sim \mathcal{GP}(0, \Phi^{\ell-1})$  and  $h^\ell$ 's explicit dynamics both depend on  $\Phi^{\ell-1}$

$$\begin{aligned}
\frac{\partial}{\partial \Phi_{\alpha\beta}^{\ell-1}(t', s')} \langle \phi(h_\mu^\ell(t)) \phi(h_\nu^\ell(s)) \rangle &= \frac{1}{2} \left\langle \frac{\partial^2}{\partial u_\alpha^\ell(t') \partial u_\beta^\ell(s')} \phi(h_\mu^\ell(t)) \phi(h_\nu^\ell(s)) \right\rangle \\
&\quad + \left\langle \frac{\partial}{\partial \Phi_{\alpha\beta}^{\ell-1}(t', s')} \phi(h_\mu^\ell(t)) \phi(h_\nu^\ell(s)) \right\rangle \quad (24)
\end{aligned}$$

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

$$\begin{aligned}
\frac{\partial^2 S}{\partial \hat{\Delta}_\mu(t) \partial \Delta_\alpha(s)} &= \delta_{\mu\alpha} \delta(t - s) + \Theta(t - s) K_{\mu\alpha}(s) \\
\frac{\partial^2}{\partial \hat{\Delta}_\mu(t) \partial K_{\alpha\beta}(s)} &= \delta_{\mu\alpha} \Theta(t - s) \Delta_\beta(s) \\
\frac{\partial^2 S}{\partial \hat{K}_{\mu\nu}(t) \partial K_{\alpha\beta}(t')} &= \delta_{\mu\alpha} \delta_{\nu\beta} \delta(t - t') \\
\frac{\partial^2 S}{\partial \hat{K}_{\mu\nu}(t) \partial \Phi_{\alpha\beta}^\ell(t', s')} &= \delta_{\mu\alpha} \delta_{\nu\beta} G_{\alpha\beta}^{\ell+1}(t', s') \delta(t - t') \delta(t - s') \\
\frac{\partial^2 S}{\partial \hat{K}_{\mu\nu}(t) \partial G_{\alpha\beta}^\ell(t', s')} &= \delta_{\mu\alpha} \delta_{\nu\beta} \Phi_{\alpha\beta}^{\ell-1}(t', s') \delta(t - t') \delta(t - s') \quad (25)
\end{aligned}$$

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

$$\mathbf{q}_1 = \text{Vec}\{\Phi_{\mu\nu}^\ell(t, s), G_{\mu\nu}^\ell(t, s), K_{\mu\nu}(t), \Delta_\mu(t), \hat{\Phi}_{\mu\nu}^\ell(t, s), \hat{G}_{\mu\nu}^\ell(t, s), \hat{K}_{\mu\nu}(t), \hat{\Delta}_\mu(t)\} \quad (26)$$

$$\mathbf{q}_2 = \text{Vec}\{A_{\mu\nu}^\ell(t, s), B_{\mu\nu}^\ell(t, s)\}. \quad (27)$$

611 This choice will become apparent shortly.

$$\nabla_{\mathbf{q}}^2 S = \begin{bmatrix} \nabla_{\mathbf{q}_1}^2 S & \nabla_{\mathbf{q}_1 \mathbf{q}_2}^2 S \\ \nabla_{\mathbf{q}_2 \mathbf{q}_1}^2 S & \nabla_{\mathbf{q}_2}^2 S \end{bmatrix} \quad (28)$$

612 To calculate the full propagator  $\Sigma = -[\nabla_{\mathbf{q}}^2 S]^{-1}$ , we will assume invertibility of the upper block

613  $\Sigma^0 = -[\nabla_{\mathbf{q}_1}^2 S]^{-1}$  and use this in the Schur complement

$$\begin{aligned} \Sigma &= -[\nabla_{\mathbf{q}}^2 S]^{-1} = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \\ \Sigma_{11} &= \Sigma^0 - \Sigma^0 [\nabla_{\mathbf{q}_1 \mathbf{q}_2}^2 S] (\nabla_{\mathbf{q}_2}^2 S + (\nabla_{\mathbf{q}_2 \mathbf{q}_1}^2 S) \Sigma^0 (\nabla_{\mathbf{q}_1 \mathbf{q}_2}^2 S))^{-1} [\nabla_{\mathbf{q}_2 \mathbf{q}_1}^2 S] \Sigma^0 \\ \Sigma_{12} &= \Sigma_{21}^\top = -\Sigma^0 [\nabla_{\mathbf{q}_1 \mathbf{q}_2}^2 S] (\nabla_{\mathbf{q}_2}^2 S + (\nabla_{\mathbf{q}_2 \mathbf{q}_1}^2 S) \Sigma^0 (\nabla_{\mathbf{q}_1 \mathbf{q}_2}^2 S))^{-1} \\ \Sigma_{22} &= -(\nabla_{\mathbf{q}_2}^2 S + (\nabla_{\mathbf{q}_2 \mathbf{q}_1}^2 S) \Sigma^0 (\nabla_{\mathbf{q}_1 \mathbf{q}_2}^2 S))^{-1} \end{aligned} \quad (29)$$

614 We now need to solve for  $\Sigma^0 = -[\nabla_{\mathbf{q}_1}^2 S]^{-1}$ . To perform this inverse, we again partition  $\mathbf{q}_1$  into

615 two sets of order parameters  $\mathbf{q}_1 = [\mathbf{q}_1^1, \mathbf{q}_1^2]$  where  $\mathbf{q}_1^1 = \text{Vec}\{\Phi_{\mu\nu}^\ell(t, s), G_{\mu\nu}^\ell(t, s), K_{\mu\nu}(t), \Delta_\mu(t)\}$

616 and  $\mathbf{q}_1^2 = \text{Vec}\{\hat{\Phi}_{\mu\nu}^\ell(t, s), \hat{G}_{\mu\nu}^\ell(t, s), \hat{K}_{\mu\nu}(t), \hat{\Delta}_\mu(t)\}$

$$\nabla_{\mathbf{q}_1}^2 S = \begin{bmatrix} \mathbf{0} & \mathbf{U}^\top \\ \mathbf{U} & \boldsymbol{\kappa} \end{bmatrix}, \quad \boldsymbol{\kappa} \equiv \nabla_{\mathbf{q}_1^2}^2 S, \quad \mathbf{U} \equiv \nabla_{\mathbf{q}_1^2 \mathbf{q}_1^1}^2 S \quad (30)$$

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

$$\Sigma^0 = -[\nabla_{\mathbf{q}_1}^2 S]^{-1} = \begin{bmatrix} \mathbf{U}^{-1} \boldsymbol{\kappa} [\mathbf{U}^{-1}]^\top & -\mathbf{U}^{-1} \\ -[\mathbf{U}^\top]^{-1} & \mathbf{0} \end{bmatrix} \quad (31)$$

619 Thus given  $\boldsymbol{\kappa}, \mathbf{U}$ , we can solve for  $\Sigma^0$  and ultimately for the full propagator  $\Sigma$ . The relevant entries

620 in  $\boldsymbol{\kappa}$  and  $\mathbf{U}$  are given by those second derivatives calculated above. We note that each of the field

621 derivatives needed for  $\mathbf{U}$  can be computed implicitly from the field dynamics. For example, for the

622  $\Delta_\mu(t)$  derivatives we have

$$\begin{aligned} \frac{\partial}{\partial \Delta_{\nu'}(t')} h_\mu^\ell(t) &= \gamma \Theta(t - t') \Phi_{\mu\nu'}^{\ell-1}(t, t') g_\nu^\ell(t') \\ &+ \gamma \int_0^t ds \sum_\nu [A_{\mu\nu}^{\ell-1}(t, s) + \Phi_{\mu\nu}^{\ell-1}(t, s) \Delta_\nu(s)] \frac{\partial g_\nu^\ell(s)}{\partial \Delta_{\nu'}(t')} \\ \frac{\partial}{\partial \Delta_\nu(t')} z_\mu^\ell(t) &= \gamma \Theta(t - t') G_{\mu\nu}^{\ell+1}(t, t') \phi(h_\nu^\ell(t')) \\ &+ \gamma \int_0^t ds \sum_\nu [B_{\mu\nu}^\ell(t, s) + G_{\mu\nu}^{\ell+1}(t, s) \Delta_\nu(s)] \frac{\partial \phi(h_\nu^\ell(s))}{\partial \Delta_{\nu'}(t')} \end{aligned} \quad (32)$$

623 These can then be used in the averages such as  $\left\langle \frac{\partial}{\partial \Delta_{\nu'}(t')} \phi(h_\mu^\ell(t)) \phi(h_\nu^\ell(s)) \right\rangle$ . Similarly, we can

624 compute terms such as  $\frac{\partial h_\mu^\ell(t)}{\partial \Phi_{\alpha\beta}^\ell(t', s')}$  through the following closed equations

$$\begin{aligned} \frac{\partial h_\mu^\ell(t)}{\partial \Phi_{\alpha\beta}^{\ell-1}(t', s')} &= \gamma \delta(t - t') \delta_{\mu\alpha} \Theta(t - s') \Delta_\beta(s') \\ &+ \gamma \int_0^t ds \sum_\nu [A_{\mu\nu}^{\ell-1}(t, s) + \Delta_\nu(s) \Phi_{\mu\nu}^{\ell-1}(t, s)] \frac{\partial g_\nu^\ell(s)}{\partial \Phi_{\alpha\beta}^\ell(t', s')} \\ \frac{\partial z_\mu^\ell(t)}{\partial \Phi_{\alpha\beta}^{\ell-1}(t', s')} &= \gamma \int_0^t ds \sum_\nu [B_{\mu\nu}^\ell(t, s) + \Delta_\nu(s) G_{\mu\nu}^{\ell+1}(t, s)] \frac{\partial \phi(h_\nu^\ell(s))}{\partial \Phi_{\alpha\beta}^{\ell-1}(t', s')} \end{aligned} \quad (33)$$

625 These terms can then be used to compute quantities like  $D^{\Phi^\ell}$ .

## 626 E Solving for the Propagator

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

---

### Algorithm 1: Propagator Solver

---

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

**Result:** Propagator Matrix  $\Sigma$

- 1 Solve DMFT equations with Algorithm 2 for order parameters  $f_\mu(t), \Phi_{\mu\alpha}^\ell(t, s), \dots$ ;
  - 2 Draw  $\mathcal{S}$  samples  $\{u_{\mu,n}^\ell(t)\}_{n=1}^{\mathcal{S}} \sim \mathcal{GP}(0, \Phi^{\ell-1}), \{r_{\mu,n}^\ell(t)\}_{n=1}^{\mathcal{S}} \sim \mathcal{GP}(0, \mathbf{G}^{\ell+1})$ ;
  - 3 Integrate dynamics for each sample to get  $\{h_{\mu,n}^\ell(t), z_{\mu,n}^\ell(t)\}_{n=1}^{\mathcal{S}}$ ;
  - 4 Estimate  $\kappa$  functions with Monte-Carlo integration, for instance
  - 5  $\kappa_{\mu\nu\alpha\beta}^{\Phi^\ell}(t, s, t', s') =$
  - 628  $\frac{1}{\mathcal{S}} \sum_{n \in [\mathcal{S}]} \phi(h_{\mu,n}^\ell(t)) \phi(h_{\nu,n}^\ell(s)) \phi(h_{\alpha,n}^\ell(t')) \phi(h_{\beta,n}^\ell(s')) - \Phi_{\mu\nu}^\ell(t, s) \Phi_{\alpha\beta}^\ell(t', s')$ ;
  - 6 For each sample, compute field sensitivities to error signals, such as  $\frac{\partial h_{\mu,n}^\ell(t)}{\partial \Delta_\nu(s)}$ , and kernels
  - $\frac{\partial h_{\mu,n}^\ell(t)}{\partial \Phi_{\alpha\beta}^\ell(t', s')}$  implicitly using equations (32) (33) ;
  - 7 Use these sensitivities to compute the necessary  $D$  tensors such as
  - $D_{\mu\nu\alpha}^{\Phi^\ell \Delta} = \frac{1}{\mathcal{S}} \sum_{n \in [\mathcal{S}]} \frac{\partial}{\partial \Delta_\alpha(t')} [\phi(h_{\mu,n}^\ell(t)) \phi(h_{\nu,n}^\ell(s))]$ ;
  - 8 Invert  $\mathbf{U}$  matrix and compute  $\Sigma_0$  in equation (31);
  - 9 Compute the Schur-complement in equation (29) to handle the response functions ;
- 

629 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  $\{\Phi^\ell, \mathbf{G}^\ell\}_{\ell=1}^L, \{\mathbf{A}^\ell, \mathbf{B}^\ell\}_{\ell=1}^{L-1}$ , Sample count  $\mathcal{S}$ , Update Speed  $\beta$

**Result:** Final Kernels  $\{\Phi^\ell, \mathbf{G}^\ell\}_{\ell=1}^L, \{\mathbf{A}^\ell, \mathbf{B}^\ell\}_{\ell=1}^{L-1}$ , Network predictions through training  $f_\mu(t)$

- 1  $\Phi^0 = K^x \otimes \mathbf{11}^\top, \mathbf{G}^{L+1} = \mathbf{11}^\top$  ;
  - 2 **while** Kernels Not Converged **do**
  - 3     From  $\{\Phi^\ell, \mathbf{G}^\ell\}$  compute  $\mathbf{K}^{NTK}(t, t)$  and solve  $\frac{d}{dt} f_\mu(t) = \sum_\alpha \Delta_\alpha(t) K_{\mu\alpha}^{NTK}(t, t)$ ;
  - 4      $\ell = 1$ ;
  - 5     **while**  $\ell < L + 1$  **do**
  - 6         Draw  $\mathcal{S}$  samples  $\{u_{\mu,n}^\ell(t)\}_{n=1}^{\mathcal{S}} \sim \mathcal{GP}(0, \Phi^{\ell-1}), \{r_{\mu,n}^\ell(t)\}_{n=1}^{\mathcal{S}} \sim \mathcal{GP}(0, \mathbf{G}^{\ell+1})$ ;
  - 7         Integrate dynamics for each sample to get  $\{h_{\mu,n}^\ell(t), z_{\mu,n}^\ell(t)\}_{n=1}^{\mathcal{S}}$ ;
  - 8         Compute new  $\Phi^\ell, \mathbf{G}^\ell$  estimates:
  - 9          $\tilde{\Phi}_{\mu\alpha}^\ell(t, s) = \frac{1}{\mathcal{S}} \sum_{n \in [\mathcal{S}]} \phi(h_{\mu,n}^\ell(t)) \phi(h_{\alpha,n}^\ell(s)), \tilde{G}_{\mu\alpha}^\ell(t, s) = \frac{1}{\mathcal{S}} \sum_{n \in [\mathcal{S}]} g_{\mu,n}^\ell(t) g_{\alpha,n}^\ell(s)$  ;
  - 10         Solve for Jacobians on each sample  $\frac{\partial \phi(h_n^\ell)}{\partial \mathbf{r}_n^{\ell\top}}, \frac{\partial g_n^\ell}{\partial \mathbf{u}_n^{\ell\top}}$  ;
  - 630 11         Compute new  $\mathbf{A}^\ell, \mathbf{B}^{\ell-1}$  estimates:
  - 12          $\tilde{\mathbf{A}}^\ell = \frac{1}{\mathcal{S}} \sum_{n \in [\mathcal{S}]} \frac{\partial \phi(h_n^\ell)}{\partial \mathbf{r}_n^{\ell\top}}, \tilde{\mathbf{B}}^{\ell-1} = \frac{1}{\mathcal{S}} \sum_{n \in [\mathcal{S}]} \frac{\partial g_n^\ell}{\partial \mathbf{u}_n^{\ell\top}}$  ;
  - 13          $\ell \leftarrow \ell + 1$ ;
  - 14     **end**
  - 15      $\ell = 1$ ;
  - 16     **while**  $\ell < L + 1$  **do**
  - 17         Update feature kernels:  $\Phi^\ell \leftarrow (1 - \beta)\Phi^\ell + \beta\tilde{\Phi}^\ell, \mathbf{G}^\ell \leftarrow (1 - \beta)\mathbf{G}^\ell + \beta\tilde{\mathbf{G}}^\ell$  ;
  - 18         **if**  $\ell < L$  **then**
  - 19             Update  $\mathbf{A}^\ell \leftarrow (1 - \beta)\mathbf{A}^\ell + \beta\tilde{\mathbf{A}}^\ell, \mathbf{B}^\ell \leftarrow (1 - \beta)\mathbf{B}^\ell + \beta\tilde{\mathbf{B}}^\ell$
  - 20         **end**
  - 21          $\ell \leftarrow \ell + 1$
  - 22     **end**
  - 23 **end**
  - 24 **return**  $\{\Phi^\ell, \mathbf{G}^\ell\}_{\ell=1}^L, \{\mathbf{A}^\ell, \mathbf{B}^\ell\}_{\ell=1}^{L-1}, \{f_\mu(t)\}_{\mu=1}^P$
-

## F 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 C)

$$\begin{aligned} \langle q_i - q_i^\infty \rangle_N &= \frac{N \langle (q_i - q_i^\infty) V \rangle_\infty + \frac{N^2}{2} \langle (q_i - q_i^\infty) V^2 \rangle_\infty \dots}{1 + N \langle V \rangle_\infty + \frac{N^2}{2} \langle V^2 \rangle_\infty + \dots} \\ &\sim \frac{1}{3!N} \sum_{jkl} \frac{\partial^3 S}{\partial q_j \partial q_k \partial q_l} \langle \delta_i \delta_j \delta_k \delta_l \rangle_\infty + \mathcal{O}(N^{-2}). \end{aligned} \quad (34)$$

$$= \frac{1}{2N} \sum_{jkl} \frac{\partial^3 S}{\partial q_j \partial q_k \partial q_l} \Sigma_{ij} \Sigma_{kl} + \mathcal{O}(N^{-2}) \quad (35)$$

where  $\delta_j = \sqrt{N}(q_j - q_j^\infty)$  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  $\Sigma_{ij}$ , which was provided in the preceding section D. In practice computing even the full set of second derivatives for the DMFT action to get  $\Sigma$  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 6 and A.6

### F.1 Correction to Mean Predictions and Full MSE Correction

Supposing that we solved for the propagator  $\Sigma$ , 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.

$$\begin{aligned} \frac{d}{dt} \langle \Delta_\mu(t) \rangle &= - \sum_\nu \langle K_{\mu\nu}(t) \Delta_\nu(t) \rangle \\ &= - \sum_\nu \langle K_{\mu\nu}(t) \rangle \langle \Delta_\nu(t) \rangle - \sum_\nu \text{Cov}(K_{\mu\nu}(t), \Delta_\nu(t)) \\ &\sim - \sum_\nu \langle K_{\mu\nu}(t) \rangle \langle \Delta_\nu(t) \rangle - \frac{1}{N} \sum_\nu \Sigma_{\mu\nu\nu}^{K\Delta}(t, t) + \mathcal{O}(N^{-2}) \end{aligned} \quad (36)$$

where  $\Sigma_{\mu\nu\nu}^{K\Delta}(t, t)$  is the leading covariance (propagator element) between the kernel  $K_{\mu\nu}(t)$  and prediction error  $\Delta_\nu(t)$ . We see that the average kernel  $\langle K_{\mu\nu}(t) \rangle$  (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  $\langle \Delta_\mu(t) \rangle$ , the square loss at width  $N$  and time  $t$  has the form

$$\sum_\mu \langle \Delta_\mu(t)^2 \rangle \sim \left(1 - \frac{2}{N}\right) \sum_\mu \Delta_\mu^\infty(t)^2 + \frac{2}{N} \sum_\mu \langle \Delta_\mu(t) \rangle_\infty \Delta_\mu^\infty(t) + \frac{1}{N} \sum_\mu \Sigma_{\mu\mu}^\Delta(t, t) + \mathcal{O}(N^{-2}) \quad (37)$$

We will now comment on the structure of the cross term in this above solution. First, if  $\langle \mathbf{K} \rangle \succeq \mathbf{K}^\infty$  and  $\Sigma^{K\Delta}$  is negligible then the average errors at finite width will decay more rapidly than the infinite width model. However, we suspect that in general,  $\langle \mathbf{K} \rangle - \mathbf{K}^\infty$  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 6 and Figure A.6. Additionally, the term involving  $\Sigma^{K\Delta}$  will generically increase the cross term since the dynamics of  $\Delta$  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 G.2).

## F.2 Perturbation Theory in Rates rather than Predictions

In experiments on deep CNNs trained on CIFAR-10 in [6] and [A.6], 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

$$\frac{d}{dt}r(t) = K(t) \quad (38)$$

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  $1 - \frac{\epsilon}{K}$ .

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

$$\frac{d}{dt}T(t) = -K(t)T(t), \quad T(0) = \mathbf{I}. \quad (39)$$

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

$$R(t) = -\log T(t) \quad (40)$$

where  $\log$  is the matrix logarithm function. Note that in general  $T(t)$  admits a Peano-Baker series solution [59-61]. In the special case where  $K(t)$  commutes with  $\bar{K}(t) = \frac{1}{t} \int_0^t ds K(s)$ , we obtain the following simplified formula for the rate matrix  $R$

$$R(t) = \int_0^t ds K(s) \quad (41)$$

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

$$\begin{aligned} \langle |\Delta(t)|^2 \rangle &= y^\top \langle \exp(-2R(t)) \rangle y \\ &\sim y^\top \exp \left( -2 \left( R_\infty(t) + \frac{1}{N} R^1(t) \right) \right) y \\ &\quad + \frac{1}{2} \sum_{\mu\nu\alpha\beta} \Sigma_{\mu\nu\alpha\beta}^R(t, t) \frac{\partial^2}{\partial R_{\mu\nu} \partial R_{\alpha\beta}} y^\top \exp(-2R) y|_{R=R_\infty(t) + \frac{1}{N} R^1(t)} + \mathcal{O}(N^{-2}) \end{aligned} \quad (42)$$

where  $\langle R \rangle \sim R_\infty + \frac{1}{N} R^1 + \mathcal{O}(N^{-2})$  is the leading correction to the mean  $R$ . In this representation, it is clear that finite width can alter the timescale of the dynamics through a correction to the mean of  $R$ , 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 [6] and [A.6]

## G 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 [G.1] which relies only on linearizing the dynamics around the infinite width solution. In the

694  $\gamma \rightarrow 0$  limit, all of the  $D$  tensors vanish and the  $\kappa$  tensors are constant in time. Thus, it suffices to  
 695 analyze the kernels restricted to  $t = 0$  and study the evolution of the prediction variance  $\Delta(t)$ .

$$\begin{aligned}
 S &= \int dt \sum_{\mu} \hat{\Delta}_{\mu}(t) \left( \Delta_{\mu}(t) - y_{\mu} + \int ds \sum_{\nu} \Theta(t-s) K_{\mu\nu} \Delta_{\nu}(s) \right) \\
 &+ \sum_{\ell} \sum_{\mu\nu} \left[ \hat{\Phi}_{\mu\nu}^{\ell} \Phi_{\mu\nu}^{\ell} + G_{\mu\nu}^{\ell} \hat{G}_{\mu\nu}^{\ell} \right] + \sum_{\mu\nu} \hat{K}_{\mu\nu} \left[ K_{\mu\nu} - \sum_{\ell} G_{\mu\nu}^{\ell+1} \Phi_{\mu\nu}^{\ell} \right] + \sum_{\ell} \ln Z_{\ell} \\
 Z_{\ell} &= \mathbb{E}_{\{u_{\mu}^{\ell}\}, \{r_{\mu}^{\ell}\}} \exp \left( - \sum_{\mu\nu} \hat{\Phi}_{\mu\nu}^{\ell} \phi(u_{\mu}^{\ell}) \phi(u_{\nu}^{\ell}) - \sum_{\mu\nu} \hat{G}_{\mu\nu}^{\ell} g_{\mu}^{\ell} g_{\nu}^{\ell} \right), \quad g_{\mu}^{\ell} = r_{\mu}^{\ell} \dot{\phi}(u_{\mu}^{\ell}) \quad (43)
 \end{aligned}$$

696 where  $\{u_{\mu}^{\ell}\} \sim \mathcal{N}(0, \Phi^{\ell-1})$ ,  $\{r_{\mu}^{\ell}\} \sim \mathcal{N}(0, G^{\ell+1})$ . Taking two derivatives with respect to  $\{\hat{\Phi}^{\ell}, \hat{G}^{\ell}\}$   
 697 give terms of the form

$$\begin{aligned}
 \kappa_{\mu\nu\alpha\beta}^{\Phi^{\ell}} &= \langle \phi(u_{\mu}^{\ell}) \phi(u_{\nu}^{\ell}) \phi(u_{\alpha}^{\ell}) \phi(u_{\beta}^{\ell}) \rangle - \Phi_{\mu\nu}^{\ell} \Phi_{\alpha\beta}^{\ell} \\
 \kappa_{\mu\nu\alpha\beta}^{G^{\ell}} &= \langle g_{\mu}^{\ell} g_{\nu}^{\ell} g_{\alpha}^{\ell} g_{\beta}^{\ell} \rangle - G_{\mu\nu}^{\ell} G_{\alpha\beta}^{\ell} \\
 \kappa_{\mu\nu\alpha\beta}^{\Phi^{\ell}, G^{\ell}} &= \langle \phi(u_{\mu}^{\ell}) \phi(u_{\nu}^{\ell}) g_{\alpha}^{\ell} g_{\beta}^{\ell} \rangle - \Phi_{\mu\nu}^{\ell} G_{\alpha\beta}^{\ell} \quad (44)
 \end{aligned}$$

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

$$\begin{aligned}
 D_{\mu\nu\alpha\beta}^{\Phi^{\ell+1}\Phi^{\ell}} &= \frac{\partial^2}{\partial \Phi_{\alpha\beta}^{\ell}} \langle \phi(u_{\mu}^{\ell+1}) \phi(u_{\nu}^{\ell+1}) \rangle, \quad D_{\mu\nu\alpha\beta}^{G^{\ell}G^{\ell+1}} = \frac{\partial}{\partial G_{\alpha\beta}^{\ell+1}} \langle g_{\mu}^{\ell} g_{\nu}^{\ell} \rangle \\
 D_{\mu\nu\alpha\beta}^{G^{\ell}\Phi^{\ell-1}} &= \frac{\partial}{\partial \Phi_{\alpha\beta}^{\ell-1}} \langle g_{\mu}^{\ell} g_{\nu}^{\ell} \rangle \quad (45)
 \end{aligned}$$

$$\begin{aligned}
 D_{\mu\nu\alpha\beta}^{K\Phi^{\ell}} &= \delta_{\mu\alpha} \delta_{\nu\beta} G_{\mu\nu}^{\ell+1}, \quad D_{\mu\nu\alpha\beta}^{KG^{\ell}} = \delta_{\mu\alpha} \delta_{\nu\beta} \Phi_{\mu\nu}^{\ell-1} \\
 D_{\mu\alpha\beta}^{\Delta K}(t) &= \int ds \Theta(t-s) \delta_{\mu\alpha} \Delta_{\beta}(s) \quad (46)
 \end{aligned}$$

699 As before we let  $\mathbf{q}_1 = \text{Vec}\{\Delta_{\mu}(t), \Phi_{\mu\nu}^{\ell}, G_{\mu\nu}^{\ell}, K_{\mu\nu}\}$  and  $\mathbf{q}_2 = \text{Vec}\{\hat{\Delta}_{\mu}(t), \hat{\Phi}_{\mu\nu}^{\ell}, \hat{G}_{\mu\nu}^{\ell}, \hat{K}_{\mu\nu}\}$ . The  
 700 propagator has the form

$$U \equiv \nabla_{\mathbf{q}_2 \mathbf{q}_1}^2 S = \begin{bmatrix} \mathbf{I} + \Theta_K & \mathbf{0} & \mathbf{0} & D^{\Delta K} \\ \mathbf{0} & \mathbf{I} - D^{\Phi, \Phi} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & -D^{G\Phi} & \mathbf{I} - D^{GG} & \mathbf{0} \\ \mathbf{0} & -D^{K\Phi} & -D^{KG} & \mathbf{I} \end{bmatrix}, \quad \nabla_{\mathbf{q}_2 \mathbf{q}_2}^2 S = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \kappa^{\Phi, \Phi} & \kappa^{\Phi G} & \mathbf{0} \\ \mathbf{0} & \kappa^{G\Phi} & \kappa^{GG} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \quad (47)$$

701 The propagator of interest is  $\Sigma_{\mathbf{q}_1} = U^{-1} [\nabla_{\mathbf{q}_2 \mathbf{q}_2}^2 S] U^{-1\top}$ . We can exploit the block structure of  $U$   
 702 to find an inverse

$$U^{-1} = \begin{bmatrix} U_{\Delta\Delta}^{-1} & U_{\Delta\Phi}^{-1} & U_{\Delta G}^{-1} & U_{\Delta K}^{-1} \\ \mathbf{0} & U_{\Phi\Phi}^{-1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & U_{G\Phi}^{-1} & U_{GG}^{-1} & \mathbf{0} \\ \mathbf{0} & U_{K\Phi}^{-1} & U_{KG}^{-1} & \mathbf{I} \end{bmatrix} \quad (48)$$

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

$$\begin{aligned}
 \Sigma &= \begin{bmatrix} \mathbf{0} & U_{\Delta\Phi}^{-1} \kappa^{\Phi\Phi} + U_{\Delta G}^{-1} \kappa^{G\Phi} & U_{\Delta\Phi}^{-1} \kappa^{\Phi G} + U_{\Delta G}^{-1} \kappa^{GG} & \mathbf{0} \\ \mathbf{0} & U_{\Phi\Phi}^{-1} \kappa^{\Phi\Phi} & U_{\Phi\Phi}^{-1} \kappa^{\Phi G} & \mathbf{0} \\ \mathbf{0} & U_{G\Phi}^{-1} \kappa^{G\Phi} + U_{GG}^{-1} \kappa^{GG} & U_{G\Phi}^{-1} \kappa^{GG} + U_{GG}^{-1} \kappa^{GG} & \mathbf{0} \\ \mathbf{0} & U_{K\Phi}^{-1} \kappa^{K\Phi} + U_{KG}^{-1} \kappa^{KG} & U_{K\Phi}^{-1} \kappa^{KG} + U_{KG}^{-1} \kappa^{GG} & \mathbf{0} \end{bmatrix} \\
 &\times \begin{bmatrix} U_{\Delta\Delta}^{-1} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ [U_{\Delta\Phi}^{-1}]^{\top} & U_{\Phi\Phi}^{-1} & [U_{G\Phi}^{-1}]^{\top} & [U_{K\Phi}^{-1}]^{\top} \\ [U_{\Delta G}^{-1}]^{\top} & \mathbf{0} & U_{GG}^{-1} & [U_{KG}^{-1}]^{\top} \\ [U_{\Delta K}^{-1}]^{\top} & \mathbf{0} & \mathbf{0} & \mathbf{I} \end{bmatrix} \quad (49)
 \end{aligned}$$



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

$$\Sigma_K = U_{K\Phi}^{-1} \kappa^{\Phi\Phi} [U_{K\Phi}^{-1}]^\top + U_{KG}^{-1} \kappa^{G\Phi} [U_{K\Phi}^{-1}]^\top + U_{K\Phi}^{-1} \kappa^{\Phi G} [U_{KG}^{-1}]^\top + U_{KG}^{-1} \kappa^{GG} [U_{KG}^{-1}]^\top \quad (50)$$

Next we use the fact that  $U_{\Delta\Phi}^{-1} = U_{\Delta K}^{-1} U_{K\Phi}^{-1}$  and that  $U_{\Delta G}^{-1} = U_{\Delta K}^{-1} U_{KG}^{-1}$ , which follows from the block structure of  $U$ . Consequently we arrive at the identity

$$\begin{aligned} \Sigma_\Delta &= U_{\Delta\Phi}^{-1} \kappa^{\Phi\Phi} [U_{\Delta\Phi}^{-1}]^{-1} + U_{\Delta G}^{-1} \kappa^{G\Phi} [U_{\Delta\Phi}^{-1}]^{-1} + U_{\Delta G}^{-1} \kappa^{\Phi G} [U_{\Delta\Phi}^{-1}]^{-1} + U_{\Delta G}^{-1} \kappa^{GG} [U_{\Delta G}^{-1}]^{-1} \\ &= U_{\Delta K}^{-1} \Sigma_K [U_{K\Delta}^{-1}]^\top. \end{aligned} \quad (51)$$

Lastly, we note that, by the Schur-complement formula that  $U_{\Delta K}^{-1} = -(\mathbf{I} + \Theta_K)^{-1} D^{\Delta K}$ . Thus, writing  $(\mathbf{I} + \Theta_K) \Sigma_\Delta (\mathbf{I} + \Theta_K)^\top = D^{\Delta K} \Sigma_K [D^{\Delta K}]^\top$  as an integral equation, we find

$$\begin{aligned} \Sigma_{\mu\nu}^\Delta(t, s) &+ \int_0^t dt' \sum_\alpha K_{\mu\alpha} \Sigma_{\alpha\nu}^\Delta(t', s) + \int_0^s ds' \sum_\beta K_{\nu\beta} \Sigma_{\mu\beta}^\Delta(t, s') \\ &+ \int_0^t dt' \int_0^s ds' \sum_{\alpha\beta} K_{\mu\alpha} K_{\nu\beta} \Sigma_{\alpha\beta}^\Delta(t', s') = \sum_{\alpha\beta} \int_0^t \Delta_\alpha(t') \int_0^s ds' \Delta_\beta(s') \Sigma_{\mu\alpha, \nu\beta}^K \end{aligned} \quad (52)$$

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

$$\begin{aligned} \frac{\partial^2}{\partial t \partial s} \Sigma_{\mu\nu}^\Delta(t, s) &+ \sum_\alpha K_{\mu\alpha} \frac{\partial}{\partial s} \Sigma_{\alpha\nu}^\Delta(t, s) + \sum_\beta K_{\nu\beta} \frac{\partial}{\partial t} \Sigma_{\mu\beta}^\Delta(t, s) \\ &+ \sum_{\alpha\beta} K_{\mu\alpha} K_{\nu\beta} \Sigma_{\alpha\beta}^\Delta(t, s) = \sum_{\alpha\beta} \Delta_\alpha(t) \Delta_\beta(s) \Sigma_{\mu\alpha, \nu\beta}^K \end{aligned} \quad (53)$$

Let  $\{\psi_k\}$  be the eigenvectors of the kernel matrix  $K$ . Projecting these dynamics on the eigenspace  $\Sigma_{k\ell}(t, s) = \psi_k^\top \Sigma(t, s) \psi_\ell$  recovers the equation in the main text

$$\left( \frac{\partial}{\partial t} + \lambda_k \right) \left( \frac{\partial}{\partial s} + \lambda_\ell \right) \Sigma_{k\ell}(t, s) = \sum_{k'\ell'} \Delta_{k'}(t) \Delta_{\ell'}(s) \Sigma_{kk', \ell\ell'}^K \quad (54)$$

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

### G.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  $\Delta(t) = \Delta_\infty(t) + \epsilon^\Delta(t)$  and  $K = K_\infty + \epsilon^K$ , and keep terms only linear in these perturbations. The perturbation  $\epsilon^K$  is fixed in time and the dynamics of  $\epsilon^\Delta(t)$  are

$$\frac{d}{dt} \epsilon^\Delta(t) = -K_\infty \epsilon^\Delta(t) - \epsilon^K \Delta_\infty(t) \quad (55)$$

Projecting this equation on the eigenspace of  $K_\infty$  gives

$$\frac{d}{dt} \epsilon_k^\Delta(t) = -\lambda_k \epsilon_k(t) - \sum_{k'} \epsilon_{kk'}^K \Delta_{k'}^\infty(t) \quad (56)$$

This immediately recovers the final result of the last section

$$\begin{aligned} N \left( \frac{\partial}{\partial t} + \lambda_k \right) \left( \frac{\partial}{\partial t} + \lambda_k \right) \langle \epsilon_k^\Delta(t) \epsilon_\ell^\Delta(s) \rangle &= \left( \frac{\partial}{\partial t} + \lambda_k \right) \left( \frac{\partial}{\partial t} + \lambda_k \right) \Sigma_{k\ell}^\Delta(t, s) \\ &= \sum_{k'\ell'} \Sigma_{kk', \ell\ell'}^K \Delta_{k'}^\infty(t) \Delta_{\ell'}^\infty(s) \end{aligned} \quad (57)$$

Qualitatively, the process of computing this linear correction (in  $\epsilon^K$ ) to the dynamics of  $\Delta$  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.

## 725 G.2 Mean Prediction Error Correction in the Lazy Limit

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

$$\begin{aligned} \frac{d}{dt} \langle \Delta(t) \rangle &= \frac{d}{dt} \Delta^\infty(t) + \frac{1}{N} \frac{d}{dt} \Delta^1(t) + \dots \\ &= -\langle (K - K^\infty + K^\infty) (\Delta - \Delta^\infty + \Delta^\infty) \rangle \\ &= -K^\infty \Delta^\infty - K^\infty \langle \Delta - \Delta^\infty \rangle \\ &\quad - \langle (K - K^\infty) \rangle \Delta^\infty - \langle (K - K^\infty) (\Delta - \Delta^\infty) \rangle \\ &\sim -K^\infty \Delta^\infty - \frac{1}{N} K^\infty \Delta^1 - \frac{1}{N} K^1 \Delta^\infty - \frac{1}{N} \langle \epsilon^K \epsilon^\Delta \rangle_\infty + \mathcal{O}(N^{-2}) \end{aligned} \quad (58)$$

728 From the previous section we have that

$$\frac{d}{dt} \epsilon^\Delta = -K^\infty \epsilon^\Delta - \epsilon^K \Delta^\infty \implies \epsilon^\Delta(t) = - \int_0^t ds \exp(-K^\infty(t-s)) \epsilon^K \exp(-K^\infty s) \mathbf{y} \quad (59)$$

729 Projecting these dynamics onto the eigenspace of the kernel gives

$$\epsilon_k^\Delta(t) = - \sum_\ell \epsilon_{k\ell}^K \frac{e^{-\lambda_\ell t} - e^{-\lambda_k t}}{\lambda_k - \lambda_\ell} y_\ell \quad (60)$$

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

$$\begin{aligned} \left( \frac{d}{dt} + \lambda_k \right) \Delta_k^1(t) &= - \sum_\ell K_{k\ell}^1 y_\ell e^{-\lambda_\ell t} + \sum_{\ell\ell'} \Sigma_{k\ell\ell'}^K \frac{e^{-\lambda_{\ell'} t} - e^{-\lambda_\ell t}}{\lambda_\ell - \lambda_{\ell'}} y_{\ell'}. \\ &= \sum_\ell y_\ell e^{-\lambda_\ell t} [-K_{k\ell}^1 + \Sigma_{k\ell\ell}^K t] + \sum_{\ell \neq \ell'} \Sigma_{k\ell\ell'}^K \frac{e^{-\lambda_{\ell'} t} - e^{-\lambda_\ell t}}{\lambda_\ell - \lambda_{\ell'}} y_{\ell'} \end{aligned} \quad (61)$$

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

$$\begin{aligned} \left( \frac{d}{dt} + \lambda \right) \Delta^1(t) &= y e^{-\lambda t} [-K^1 + t \Sigma^K] \implies \Delta^1(t) = y \left[ -t K^1 + \frac{1}{2} t^2 \Sigma^K \right] e^{-\lambda t} \\ \implies \langle \Delta(t)^2 \rangle &\sim \Delta^\infty(t)^2 + \frac{1}{N} \left[ 2y^2 t e^{-2\lambda t} \left[ -K^1 + \frac{1}{2} t \Sigma^K \right] + \Sigma^\Delta(t, t) \right] + \mathcal{O}(N^{-2}) \\ &\sim \Delta^\infty(t)^2 + \frac{2}{N} y^2 t e^{-2\lambda t} [-K^1 + \Sigma^K t] + \mathcal{O}(N^{-2}) \end{aligned} \quad (62)$$

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

## 739 H Two Layer Equations and Time/Time Diagonal

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

## 744 H.1 A Single Training Point

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

$$\begin{aligned}
 S[\{K(t), \hat{K}(t), \Delta(t), \hat{\Delta}(t)\}] & \quad (63) \\
 &= \int dt \left[ K(t) \hat{K}(t) + \hat{\Delta}(t) \left( \Delta(t) - y + \int ds \Theta(t-s) \Delta(s) K(s) \right) \right] \\
 &\quad + \ln \mathcal{Z}[\hat{K}, f], \quad \mathcal{Z} = \mathbb{E}_{h,g} \exp \left( - \int dt \hat{K}(t) [\phi(h(t))^2 + g(t)^2] \right).
 \end{aligned}$$

747 The saddle point equations are

$$\begin{aligned}
 \frac{\partial S}{\partial \hat{K}(t)} &= K(t) - \langle [\phi(h(t))^2 + g(t)^2] \rangle = 0 \\
 \frac{\partial S}{\partial \hat{\Delta}(t)} &= \Delta(t) - y + \int ds \Theta(t-s) \Delta(s) K(s) = 0 \\
 \frac{\partial S}{\partial K(s)} &= \hat{K}(s) + \Delta(s) \int dt \hat{\Delta}(t) \Theta(t-s) = 0 \\
 \frac{\partial S}{\partial \Delta(s)} &= \hat{\Delta}(s) + K(s) \int dt \hat{\Delta}(t) \Theta(t-s) = 0
 \end{aligned} \quad (64)$$

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

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

$$\begin{aligned}
 \frac{\partial^2 S}{\partial \mathbf{q}(t) \partial \mathbf{q}(s)^\top} &= \mathbf{0} \\
 \frac{\partial^2 S}{\partial \hat{\mathbf{q}}(t) \partial \mathbf{q}(s)^\top} &= \begin{bmatrix} \delta(t-s) + \Theta(t-s) K(s) & \Theta(t-s) \Delta(s) \\ -\left\langle \frac{\partial}{\partial \Delta(s)} (\phi(h(t))^2 + g(t)^2) \right\rangle & \delta(t-s) \end{bmatrix} \\
 \frac{\partial^2 S}{\partial \hat{\mathbf{q}}(t) \partial \hat{\mathbf{q}}(s)^\top} &= \begin{bmatrix} 0 & 0 \\ 0 & \kappa(t, s) \end{bmatrix}
 \end{aligned} \quad (65)$$

750 where  $\kappa(t, s) = \langle (\phi(h(t))^2 + g(t)^2)(\phi(h(s))^2 + g(s)^2) \rangle - K(t)K(s)$  is the NTK's fourth cumulant.

751 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

752 express the full Hessian

$$\nabla^2 S = \begin{bmatrix} \mathbf{0} & \frac{\partial^2 S}{\partial \mathbf{q} \partial \hat{\mathbf{q}}^\top} \\ \frac{\partial^2 S}{\partial \hat{\mathbf{q}} \partial \mathbf{q}^\top} & \frac{\partial^2 S}{\partial \hat{\mathbf{q}} \partial \hat{\mathbf{q}}^\top} \end{bmatrix} \implies -[\nabla^2 S]^{-1} = \begin{bmatrix} \left(\frac{\partial^2 S}{\partial \hat{\mathbf{q}} \partial \mathbf{q}^\top}\right)^{-1} \frac{\partial^2 S}{\partial \hat{\mathbf{q}} \partial \hat{\mathbf{q}}^\top} \left(\frac{\partial^2 S}{\partial \mathbf{q} \partial \hat{\mathbf{q}}^\top}\right)^{-1} & -\left(\frac{\partial^2 S}{\partial \hat{\mathbf{q}} \partial \mathbf{q}^\top}\right)^{-1} \\ -\left(\frac{\partial^2 S}{\partial \mathbf{q} \partial \hat{\mathbf{q}}^\top}\right)^{-1} & \mathbf{0} \end{bmatrix} \quad (66)$$

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

$$\Sigma_{\mathbf{q}} = \begin{bmatrix} \mathbf{I} + \Theta_K & \Theta_\Delta \\ -D & \mathbf{I} \end{bmatrix}^{-1} \begin{bmatrix} 0 & 0 \\ 0 & \kappa \end{bmatrix} \begin{bmatrix} \mathbf{I} + \Theta_K & \Theta_\Delta \\ -D & \mathbf{I} \end{bmatrix}^{-1\top}. \quad (67)$$

754 where  $[\Theta_K](t, s) = \Theta(t-s)K(s)$  and  $[\Theta_\Delta](t, s) = \Theta(t-s)\Delta(s)$ . The above equations allow  
 755 one to use the infinite width DMFT dynamics for  $K(t), \Delta(t)$  to compute the finite size fluctuation  
 756 dynamics of the kernel  $K$  and the error signal  $\Delta$ .

### 757 H.1.1 Computing Field Sensitivities

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

$$h(t) = u + \gamma \int_0^t ds \Delta(s) g(s), \quad z(t) = r + \gamma \int_0^t ds \Delta(s) \phi(h(t)). \quad (68)$$

760 Now, differentiating both sides with respect to  $\Delta(s')$  gives

$$\begin{aligned}\frac{\partial h(t)}{\partial \Delta(s')} &= \gamma \Theta(t-s') g(s') + \gamma \int_0^t ds \Delta(s) \frac{\partial g(s)}{\partial \Delta(s')} \\ \frac{\partial z(t)}{\partial \Delta(s')} &= \gamma \Theta(t-s') \phi(h(s')) + \gamma \int_0^t ds \Delta(s) \frac{\partial \phi(h(s))}{\partial \Delta(s')}.\end{aligned}\quad (69)$$

761 We can compute  $D$  Monte carlo by iteratively solving the above equations for each sampled trajectory  
762  $\{h(t), z(t)\}$  [62, 46]. Averaging the necessary fields over the Monte-carlo samples will give us the  
763 final expressions for  $D(t, s)$ .

$$D(t, s) = \left\langle \frac{\partial}{\partial \Delta(s)} (\phi(h(t))^2 + g(t)^2) \right\rangle \quad (70)$$

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

## 766 H.2 Test Point Fluctuation Dynamics

767 We now are in a position to calculate the test/train kernel and test prediction fluctuations. To do this  
768 systematically, we augment  $S$  with the test point prediction  $f_*$  and field  $h_*$  and introduce the kernel  
769  $K_*(t) = \langle \phi(h(t)) \phi(h_*(t)) + g(t) g_*(t) \rangle$ . The test prediction  $f_*$  and field  $h_*$  have dynamics

$$\begin{aligned}h_*(t) &= u_* + \gamma \int_0^t ds \Delta(s) \dot{\phi}(h_*(s)) z(s) K_*^x, \quad \langle u_* u \rangle = K_*^x \\ \frac{\partial}{\partial t} f_*(t) &= K_*(t) \Delta(t), \quad K_*(t) = \langle \phi(h(t)) \phi(h_*(t)) + g(t) g_*(t) \rangle\end{aligned}\quad (71)$$

770 The augmented action for this DMFT has the form

$$\begin{aligned}S &= \int dt \hat{f}_*(t) \left( f_*(t) - \int ds \Theta(t-s) \Delta(s) K_*(s) \right) + \int dt \hat{K}_*(t) K_*(t) \\ &+ \int dt \hat{\Delta}(t) \left( \Delta(t) - y + \int ds \Theta(t-s) \Delta(s) K(s) \right) + \int dt \hat{K}(t) K(t) \\ &+ \ln \mathbb{E} \exp \left( - \int \hat{K}(t) (\phi(h(t))^2 + g(t)^2) - \int \hat{K}_*(t) (\phi(h(t)) \phi(h_*(t)) + g(t) g_*(t)) \right)\end{aligned}\quad (72)$$

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

$$\begin{aligned}\nabla_{\hat{\mathbf{q}}\hat{\mathbf{q}}}^2 S[\mathbf{q}, \hat{\mathbf{q}}] &= \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & \boldsymbol{\kappa} & \boldsymbol{\kappa}_*^\top \\ 0 & 0 & \boldsymbol{\kappa}_* & \boldsymbol{\kappa}_{**} \end{bmatrix}, \quad \nabla_{\hat{\mathbf{q}}, \mathbf{q}}^2 S[\mathbf{q}, \hat{\mathbf{q}}] = \begin{bmatrix} \mathbf{I} + \boldsymbol{\Theta}_K & 0 & \boldsymbol{\Theta}_\Delta & 0 \\ -\boldsymbol{\Theta}_{K_*} & \mathbf{I} & 0 & -\boldsymbol{\Theta}_\Delta \\ -\mathbf{D} & 0 & \mathbf{I} & 0 \\ -\mathbf{D}_* & 0 & 0 & \mathbf{I} \end{bmatrix} \\ D(t, s) &= \left\langle \frac{\partial}{\partial \Delta(s)} (\phi(h(t))^2 + g(t)^2) \right\rangle\end{aligned}\quad (73)$$

$$D_*(t, s) = \left\langle \frac{\partial}{\partial \Delta(s)} (\phi(h(t)) \phi(h_*(t)) + g(t) g_*(t)) \right\rangle \quad (74)$$

772 Our total covariance matrix / propagator is thus

$$\Sigma = \begin{bmatrix} \mathbf{I} + \boldsymbol{\Theta}_K & 0 & \boldsymbol{\Theta}_\Delta & 0 \\ -\boldsymbol{\Theta}_{K_*} & \mathbf{I} & 0 & -\boldsymbol{\Theta}_\Delta \\ -\mathbf{D} & 0 & \mathbf{I} & 0 \\ -\mathbf{D}_* & 0 & 0 & \mathbf{I} \end{bmatrix}^{-1} \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & \boldsymbol{\kappa} & \boldsymbol{\kappa}_*^\top \\ 0 & 0 & \boldsymbol{\kappa}_* & \boldsymbol{\kappa}_{**} \end{bmatrix} \begin{bmatrix} \mathbf{I} + \boldsymbol{\Theta}_K & 0 & \boldsymbol{\Theta}_\Delta & 0 \\ -\boldsymbol{\Theta}_{K_*} & \mathbf{I} & 0 & -\boldsymbol{\Theta}_\Delta \\ -\mathbf{D} & 0 & \mathbf{I} & 0 \\ -\mathbf{D}_* & 0 & 0 & \mathbf{I} \end{bmatrix}^{-1\top} \quad (75)$$

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

### 774 H.3 Two Layer Linear Network Closed Form

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

$$\frac{dh(t)}{dt} = \gamma \Delta(t) z(t), \quad \frac{dz(t)}{dt} = \gamma \Delta(t) h(t) \quad (76)$$

777 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  
778 note that  $v_+(0) = \frac{1}{\sqrt{2}}(u + r)$  and  $v_-(0) = \frac{1}{\sqrt{2}}(u - r)$  are independent Gaussians. These functions  
779  $v_+(t), v_-(t)$  satisfy dynamics

$$\begin{aligned} \frac{dv_+}{dt} &= \gamma \Delta(t) v_+(t), \quad \frac{dv_-}{dt} = -\gamma \Delta(t) v_-(t) \\ \implies v_+(t) &= \exp\left(\gamma \int_0^t ds \Delta(s)\right) v_+(0) \implies \frac{\partial v_+(t)}{\partial \Delta(s)} = \gamma v_+(t) \Theta(t-s) \\ \implies v_-(t) &= \exp\left(-\gamma \int_0^t ds \Delta(s)\right) v_-(0) \implies \frac{\partial v_-(t)}{\partial \Delta(s)} = -\gamma v_-(t) \Theta(t-s) \end{aligned} \quad (77)$$

780 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  
781 normal random variables to compute  $K(t) = \langle h(t)^2 + z(t)^2 \rangle = \langle v_+(t)^2 + v_-(t)^2 \rangle$

$$\begin{aligned} D(t, s) &= \frac{\partial}{\partial \Delta(s)} \langle h(t)^2 + z(t)^2 \rangle = 2\gamma [\langle v_+(t)^2 \rangle - \langle v_-(t)^2 \rangle] \Theta(t-s) \\ &= 2\gamma \left[ \exp\left(2\gamma \int_0^t ds \Delta(s)\right) - \exp\left(-2\gamma \int_0^t ds \Delta(s)\right) \right] \Theta(t-s) \end{aligned} \quad (78)$$

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

$$\begin{aligned} \kappa(t, s) &= \langle (h(t)^2 + z(t)^2)(h(s)^2 + z(s)^2) \rangle - K(t)K(s) \\ &= 2 \langle h(t)h(s) \rangle^2 + 2 \langle h(t)z(s) \rangle^2 + 2 \langle z(t)h(s) \rangle^2 + 2 \langle z(t)z(s) \rangle^2 \\ &= \langle v_+(t)v_+(s) + v_-(t)v_-(s) \rangle^2 + \langle v_+(t)v_-(s) - v_-(t)v_+(s) \rangle^2 \end{aligned} \quad (79)$$

786 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  
787 to easily compute the single site average above.

## 788 I Multiple Samples with Whitened Data

789 In this section, we analyze the role that sample number plays in dynamics in a simplified model of a  
790 two layer linear network trained on whitened data. Concretely, we assume that  $\frac{\mathbf{x}_\mu \cdot \mathbf{x}_\nu}{D} = \delta_{\mu\nu}$ . The  
791 field equations for preactivations  $h_\mu(t)$  and pregradients  $z(t)$  obey

$$\frac{d}{dt} h_\mu(t) = \gamma \Delta_\mu(t) z(t), \quad \frac{d}{dt} z(t) = \gamma \sum_{\mu=1}^P \Delta_\mu(t) h_\mu(t) \quad (80)$$

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

$$\frac{d}{dt} h_y(t) = \Delta_y(t) z(t), \quad \frac{d}{dt} z(t) = \Delta_y(t) h_y(t), \quad \Delta_\perp(t) = 0, \quad h_\perp \sim \mathcal{N}(0, 1) \quad (81)$$

796 However, at finite width  $N$ , the off-target predictions  $\Delta_\perp$  fluctuate over random initialization. To  
797 model all of the fluctuations simultaneously, we consider the following action

$$S = \gamma \int dt \sum_{\mu} \hat{\Delta}_\mu(t) (\Delta_\mu(t) - y_\mu) + \ln \mathbb{E} \exp \left( \int dt \sum_{\mu} \hat{\Delta}_\mu(t) z(t) h_\mu(t) \right) \quad (82)$$

798 which enforces the constraint that  $\Delta_\mu(t) = y_\mu - \frac{1}{\gamma} \langle z(t) h_\mu(t) \rangle$  at infinite width. The Hessian over  
 799 order parameters  $\mathbf{q} = \text{Vec}\{\Delta_\mu(t), \hat{\Delta}_\mu(t)\}$  has the form

$$\nabla_{\mathbf{q}}^2 S = \begin{bmatrix} \mathbf{0} & (\gamma \mathbf{I} + \mathbf{D})^\top \\ \gamma \mathbf{I} + \mathbf{D} & \boldsymbol{\kappa} \end{bmatrix}, \quad D_{\mu\nu}(t, s) = \left\langle \frac{\partial}{\partial \Delta_\nu(s)} z(t) h_\mu(t) \right\rangle \quad (83)$$

800 We thus get the following covariance for predictions  $\boldsymbol{\Sigma}_\Delta = (\gamma \mathbf{I} + \mathbf{D})^{-1} \boldsymbol{\kappa} [(\gamma \mathbf{I} + \mathbf{D})^{-1}]^\top$ . We now  
 801 compute the necessary components of the  $D$  tensor

$$\begin{aligned} \frac{\partial h_\mu(t)}{\partial \Delta_\nu(s)} &= \gamma \delta_{\mu\nu} \Theta(t-s) z(s) + \gamma \int_0^t dt' \Delta_\mu(t') \frac{\partial z(t')}{\partial \Delta_\nu(s)} \\ \frac{\partial z(t)}{\partial \Delta_\nu(s)} &= \gamma \Theta(t-s) h_\nu(s) + \gamma \int_0^t dt' \sum_\mu \Delta_\mu(t') \frac{\partial h_\mu(t')}{\partial \Delta_\nu(s)} \\ &= \gamma \Theta(t-s) h_\nu(s) + \gamma \int_0^t dt' \Delta_y(t') \frac{\partial h_y(t')}{\partial \Delta_\nu(s)} \end{aligned} \quad (84)$$

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

$$\begin{aligned} \left\langle \frac{\partial h_\mu(t)}{\partial \Delta_\nu(s)} z(t) \right\rangle &= \delta_{\mu\nu} \gamma \Theta(t-s) \langle z(s) z(t) \rangle, \quad \mu, \nu \neq y \\ \left\langle \frac{\partial z(t)}{\partial \Delta_\nu(s)} h_\mu(t) \right\rangle &= \gamma \Theta(t-s) \delta_{\mu\nu}, \quad \mu, \nu \neq y \\ \left\langle \frac{\partial h_y(t)}{\partial \Delta_y(s)} z(t) \right\rangle &= \gamma \Theta(t-s) \langle z(s) z(t) \rangle + \gamma \int_0^t dt' \Delta_y(t') \left\langle \frac{\partial z(t')}{\partial \Delta_y(s)} z(t) \right\rangle \\ \left\langle \frac{\partial z(t)}{\partial \Delta_y(s)} z(t') \right\rangle &= \gamma \Theta(t-s) \langle h_y(s) z(t) \rangle + \gamma \int_0^t dt'' \Delta_y(t'') \left\langle \frac{\partial h_y(t'')}{\partial \Delta_y(s)} z(t') \right\rangle \end{aligned} \quad (85)$$

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

$$D_y(t, s) = \left\langle \frac{\partial h_y(t)}{\partial \Delta_y(s)} z(t) \right\rangle + \left\langle \frac{\partial z(t)}{\partial \Delta_y(s)} h_y(t) \right\rangle \quad (86)$$

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

$$\begin{aligned} D_\perp(t, s) &= \left\langle \frac{\partial h_\perp(t)}{\partial \Delta_\perp(s)} z(t) \right\rangle + \left\langle \frac{\partial z(t)}{\partial \Delta_\perp(s)} h_\perp(t) \right\rangle \\ &= \gamma \Theta(t-s) \langle z(t) z(s) \rangle + \gamma \Theta(t-s) \end{aligned} \quad (87)$$

810 Now, we return to  $D_y$ . To solve these equations we utilize the change of variables employed in the  
 811 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 H.3). This  
 812 orthogonal transformation decouples the dynamics

$$\frac{d}{dt} v_+(t) = \gamma \Delta_y(t) v_+(t), \quad \frac{d}{dt} v_-(t) = -\gamma \Delta_y(t) v_-(t) \quad (88)$$

813 As a consequence, the field derivatives close

$$\begin{aligned} \frac{\partial v_+(t)}{\partial \Delta_y(s)} &= \gamma \Theta(t-s) v_+(s) + \int_0^t dt' \Delta_y(t') \frac{\partial v_+(t')}{\partial \Delta_y(s)} \\ \frac{\partial v_-(t)}{\partial \Delta_y(s)} &= -\gamma \Theta(t-s) v_-(s) - \int_0^t dt' \Delta_y(t') \frac{\partial v_-(t')}{\partial \Delta_y(s)} \end{aligned} \quad (89)$$

814 The correlator of interest is

$$\langle h_y(t)z(t) \rangle = \frac{1}{2} \langle [v_+(t) + v_-(t)][v_+(t) - v_-(t)] \rangle = \frac{1}{2} \langle v_+(t)^2 - v_-(t)^2 \rangle \quad (90)$$

815 So we get that

$$\begin{aligned} D_y(t, s) &= \frac{1}{2} \left\langle \frac{\partial}{\partial \Delta_y(s)} (v_+(t)^2 - v_-(t)^2) \right\rangle \\ &= \left\langle v_+(t) \frac{\partial v_+(t)}{\partial \Delta_y(s)} \right\rangle - \left\langle v_-(t) \frac{\partial v_-(t)}{\partial \Delta_y(s)} \right\rangle \end{aligned} \quad (91)$$

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

$$\begin{aligned} \kappa_y(t, s) &= \langle v_+(t)v_+(s) + v_-(t)v_-(s) \rangle^2 + \langle v_+(t)v_+(s) - v_-(t)v_-(s) \rangle^2 \\ \kappa_\perp(t, s) &= \frac{1}{2} \langle v_+(t)v_+(s) + v_-(t)v_-(s) \rangle \end{aligned} \quad (92)$$

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

$$\begin{aligned} \Sigma_{\Delta_y} &= (\gamma \mathbf{I} + \mathbf{D}_y)^{-1} \kappa_y (\gamma \mathbf{I} + \mathbf{D}_y)^{-1} \\ \Sigma_{\Delta_\perp} &= (\gamma \mathbf{I} + \mathbf{D}_\perp)^{-1} \kappa_\perp (\gamma \mathbf{I} + \mathbf{D}_\perp)^{-1} \end{aligned} \quad (93)$$

819 Using the fact that all  $\Delta_\perp$  variables are independent and identically distributed under the leading  
820 order picture, the expected training loss has the form

$$\langle |\Delta|^2 \rangle \approx \Delta_y^\infty(t)^2 + \frac{2}{N} \Delta_y^1(t) \Delta_y^\infty(t) + \frac{1}{N} \Sigma_{\Delta_y}(t, t) + \frac{(P-1)}{N} \Sigma_{\Delta_\perp}(t, t) + \mathcal{O}(N^{-2}). \quad (94)$$

821 where  $\langle \Delta_y - \Delta_y^\infty \rangle = \frac{1}{N} \Delta_y^1(t) + \mathcal{O}(N^{-2})$ . We note that the bias correction is  $\mathcal{O}(N^{-1})$  while the  
822 variance is  $\mathcal{O}(P/N)$ . We compare the above leading order theory with and without the bias correction  
823 in Appendix Figure [A.4](#).

## 824 J Online Learning

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

$$\frac{d}{dt} \Delta(\mathbf{x}, t) = -\mathbb{E}_{\mathbf{x}' \sim p(\mathbf{x}')} K(\mathbf{x}, \mathbf{x}'; t) \Delta(\mathbf{x}', t) \quad (95)$$

829 where  $K(\mathbf{x}, \mathbf{x}'; t)$  is the dynamic NTK and  $\Delta(\mathbf{x}, t) = y(\mathbf{x}) - f(\mathbf{x}, t)$  is the prediction error. In  
830 general the distribution involves integration over an uncountable set of possible inputs  $\mathbf{x}$ . To remedy  
831 this, we utilize a countable orthonormal basis of functions for the data distribution  $\{\psi_k(\mathbf{x})\}_{k=1}^\infty$ .  
832 For example, if  $p(\mathbf{x})$  were the isotropic Gaussian density for  $\mathcal{N}(0, \mathbf{I})$ , then  $\psi_k$  could be Hermite  
833 polynomials. We expand  $\Delta$  and  $K$  in this basis  $\psi_k$ , and arrive at the following differential equation

$$\frac{d}{dt} \Delta_k(t) = - \sum_\ell K_{k\ell}(t) \Delta_\ell(t) \quad (96)$$

834 By orthonormality, the average turned into a sum over all possible orthonormal functions  $\{\psi_k\}$ .  
835 We note that since  $K$  is evolving in time, there is not generally a fixed basis of functions that  
836 diagonalize  $K$ , resulting in the couplings across eigenmodes in Equation [\(96\)](#). Since, in online  
837 learning, there is no distinction between the training and test distribution, our error of interest is  
838 simply  $\mathcal{L}(t) = \sum_k \Delta_k(t)^2$ . To obtain the finite size corrections to this quantity, we compute the joint  
839 propagator for all variables  $\{K_{k\ell}(t), \Delta_k(t)\}$ . If we wanted to pursue a perturbation theory in rates  
840 (Appendix [F.2](#)), we could again define a transition matrix  $\mathbf{T}$  and rate matrix  $\mathbf{R}(t)$  as

$$\mathbf{R}(t) = -\log \mathbf{T}(t), \quad \frac{d}{dt} T_{k\ell}(t) = - \sum_{k'} K_{kk'}(t) T_{k'\ell}(t), \quad T_{k\ell}(0) = \delta_{k\ell} \quad (97)$$

841 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  
842 correction and finite size fluctuations, so too does the error  $\Delta_k(t)$  and the loss  $\mathcal{L}$  (Appendix [F.2](#)).

## 843 J.1 Two Layer Networks

844 In the two layer case, instead of tracking kernels, we could instead deal with the distribution over  
 845 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  
 846 [6, 63]. When training on the population risk equations for  $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I})$

$$\begin{aligned} \frac{d}{dt} \mathbf{w} &= a \mathbb{E}_{\mathbf{x}} \Delta(\mathbf{x}) \dot{\phi}(\mathbf{w} \cdot \mathbf{x}) \mathbf{x} = \mathbb{E}_{\mathbf{x}} \frac{\partial \Delta(\mathbf{x})}{\partial \mathbf{x}} \dot{\phi}(\mathbf{w} \cdot \mathbf{x}) + \mathbb{E} \Delta(\mathbf{x}) \ddot{\phi}(\mathbf{w} \cdot \mathbf{x}) \mathbf{w} \\ \frac{d}{dt} a &= \mathbb{E}_{\mathbf{x}} \Delta(\mathbf{x}) \phi(\mathbf{w} \cdot \mathbf{x}) \end{aligned} \quad (98)$$

847 The action has the form

$$S = \gamma \int dt d\mathbf{x} \hat{\Delta}(t, \mathbf{x}) (\Delta(t, \mathbf{x}) - y(\mathbf{x})) + \ln \mathbb{E}_{a, \mathbf{w}} \exp \left( \int dt d\mathbf{x} \hat{\Delta}(t, \mathbf{x}) a(t) \phi(\mathbf{w}(t) \cdot \mathbf{x}) \right) \quad (99)$$

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

$$\nabla^2 S = \begin{bmatrix} 0 & \mathbf{I} + \mathbf{D}_{\Delta} \\ \mathbf{I} + \mathbf{D}_{\Delta} & \boldsymbol{\kappa} \end{bmatrix}. \quad (100)$$

849 where  $D_{\Delta}(t, \mathbf{x}; s, \mathbf{x}') = \left\langle \frac{\partial}{\partial \Delta(s, \mathbf{x}')} a(t) \phi(\mathbf{w}(t) \cdot \mathbf{x}) \right\rangle$  We can use the following implicit rule

$$\begin{aligned} \frac{\partial a(t)}{\partial \Delta(s, \mathbf{x})} &= \gamma \Theta(t - s) p(\mathbf{x}) \phi(\mathbf{w}(s) \cdot \mathbf{x}) + \gamma \mathbb{E}_{\mathbf{x}'} \int_0^t dt' \Delta(t', \mathbf{x}') \dot{\phi}(\mathbf{w} \cdot \mathbf{x}') \mathbf{x}' \cdot \frac{\partial \mathbf{w}(t)}{\partial \Delta(s, \mathbf{x})} \\ \frac{\partial \mathbf{w}(t)}{\partial \Delta(s, \mathbf{x})} &= \gamma \Theta(t - s) p(\mathbf{x}) a(s) \dot{\phi}(\mathbf{w}(s) \cdot \mathbf{x}) \mathbf{x} \\ &\quad + \gamma \mathbb{E}_{\mathbf{x}'} \int_0^t dt' \Delta(t', \mathbf{x}') \left[ \frac{\partial a(t')}{\partial \Delta(s, \mathbf{x})} \dot{\phi}(\mathbf{w} \cdot \mathbf{x}') + a(t') \ddot{\phi}(\mathbf{w} \cdot \mathbf{x}') \frac{\partial \mathbf{w}(t')}{\partial \Delta(s, \mathbf{x})} \cdot \mathbf{x}' \right] \end{aligned} \quad (101)$$

850 The above equations could be solved and then used to compute  $D_{\Delta}(t, \mathbf{x}; s, \mathbf{x}')$  which must then be  
 851 inverted to get the observed prediction variance.

## 852 J.2 Linear Activations

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

$$\begin{aligned} \frac{d}{dt} a(t) &= \gamma (\beta_{\star} - \beta(t)) \cdot \mathbf{w}(t) \\ \frac{d}{dt} \mathbf{w}(t) &= \gamma a(t) (\beta_{\star} - \beta(t)) \\ \beta(t) &= \frac{1}{\gamma} \langle a(t) \mathbf{w}(t) \rangle \end{aligned} \quad (102)$$

857 At infinite width, we see that the dynamics can be reduced to tracking the projection of the weights  
 858  $\mathbf{w}$  and  $\beta$  on the  $\beta_{\star}$  direction. The  $D - 1$  off-target dimensions vanish  $\beta_{\perp}(t) = 0$ . At infinite width,  
 859 we arrive at the alignment dynamics studied in prior work [61, 9]

$$\begin{aligned} \frac{d}{dt} \beta(t) &= \mathbf{M}(t) (\beta_{\star} - \beta(t)) \\ \frac{d}{dt} \mathbf{M}(t) &= \gamma^2 \beta(t) (\beta_{\star} - \beta(t))^{\top} + \gamma^2 \beta(t) (\beta_{\star} - \beta(t))^{\top} \\ &\quad + 2\gamma^2 (\beta_{\star} - \beta(t)) \cdot \beta(t) \mathbf{I} \end{aligned} \quad (103)$$

860 We note that  $\beta(t) = \beta(t) \beta_{\star}$  and that  $\mathbf{M}$  has only one special eigenvector  $\beta_{\star}$  with eigenvalue  $m_{\star}(t)$ .  
 861 It thus suffices to track evolution in this single direction

$$\frac{d}{dt} \beta(t) = m_{\star}(t) (\beta_{\star} - \beta(t)), \quad \frac{d}{dt} m_{\star}(t) = 4\gamma^2 \beta(t) (\beta_{\star} - \beta(t)) \quad (104)$$



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

$$\frac{d}{dt} \beta(t) = 2\sqrt{1 + \gamma^2 \beta(t)^2} (\beta_\star - \beta(t)) \quad (105)$$

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

$$f(t) \leftrightarrow \beta(t), \quad K(t) \leftrightarrow m_\star(t) \quad (106)$$

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

$$S = \gamma \int dt \hat{\beta}(t) \cdot \beta(t) + \ln \mathbb{E} \exp \left( - \int dt \hat{\beta}(t) \cdot \mathbf{w}(t) a(t) \right) \quad (107)$$

The necessary ingredients are

$$\begin{aligned} \kappa(t, s) &= \langle a(t) a(s) \mathbf{w}(t) \mathbf{w}(s)^\top \rangle - \gamma^2 \beta(t) \beta(s) \\ &= \langle a(t) a(s) \rangle \langle \mathbf{w}(t) \mathbf{w}(s)^\top \rangle + \langle a(s) \mathbf{w}(t) \rangle \langle a(t) \mathbf{w}(s)^\top \rangle \in \mathbb{R}^{D \times D} \end{aligned} \quad (108)$$

Similarly we have to compute the sensitivity tensor

$$\mathbf{D}(t, s) = \left\langle \frac{\partial}{\partial \beta(s)^\top} a(t) \mathbf{w}(t) \right\rangle \in \mathbb{R}^{D \times D} \quad (109)$$

We start from the dynamics

$$\frac{d}{dt} \mathbf{w}(t) = \gamma a(t) (\beta_\star - \beta(t)), \quad \frac{d}{dt} a(t) = \gamma (\beta_\star - \beta(t)) \cdot \mathbf{w}(t) \quad (110)$$

Next, we have to calculate causal derivatives for fields

$$\begin{aligned} \frac{\partial}{\partial \beta(s)^\top} \mathbf{w}(t) &= -\gamma \Theta(t - s) a(s) \mathbf{I} + \gamma \int_0^t dt' (\beta_\star - \beta(t')) \frac{\partial a(t')}{\partial \beta(s)^\top} \\ \frac{\partial}{\partial \beta(s)} a(t) &= -\gamma \Theta(t - s) \mathbf{w}(s) + \gamma \int_0^t dt' (\beta_\star - \beta(t')) \cdot \frac{\partial \mathbf{w}(t')}{\partial \beta(s)} \end{aligned} \quad (111)$$

Following an identical argument as in [I](#), we see that  $\mathbf{D}$  has block diagonal structure with  $D_{\beta_\star}(t, s)$  on the  $\beta_\star \beta_\star^\top$  direction and  $D_\perp(t, s)$  in any of the  $D - 1$  remaining directions

$$D_{\beta_\star}(t, s) = \left\langle \frac{\partial}{\partial \beta(s)} a(t) w_{\beta_\star}(t) \right\rangle, \quad D_\perp(t, s) = \left\langle \frac{\partial}{\partial \beta_\perp(s)} a(t) w_\perp(t) \right\rangle \quad (112)$$

Similarly,  $\kappa(t, s)$  has a similar decomposition

$$\begin{aligned} \kappa_{\beta_\star}(t, s) &= \langle a(t) a(s) \rangle \langle w_{\beta_\star}(t) w_{\beta_\star}(s) \rangle + \langle a(s) w_{\beta_\star}(t) \rangle \langle a(t) w_{\beta_\star}(s) \rangle \\ \kappa_\perp(t, s) &= \langle a(t) a(s) \rangle \langle w_\perp(t) w_\perp(s) \rangle + \langle a(s) w_\perp(t) \rangle \langle a(t) w_\perp(s) \rangle \end{aligned} \quad (113)$$

The processes have the following equations at infinite width

$$\frac{d}{dt} w_{\beta_\star}(t) = \gamma a(t) (\beta_\star - \beta(t)), \quad \frac{d}{dt} a(t) = \gamma w_{\beta_\star}(t) (\beta_\star - \beta(t)), \quad \frac{d}{dt} w_\perp(t) = 0 \quad (114)$$

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

$$\frac{d}{dt} v_+(t) = \gamma (\beta_\star - \beta(t)) v_+(t), \quad \frac{d}{dt} v_-(t) = -\gamma (\beta_\star - \beta(t)) v_-(t) \quad (115)$$

880 We can use these equations to perform the necessary averages for  $\kappa_{\beta_\star}$  and  $D_{\beta_\star}$ . Lastly, we use

$$\frac{\partial}{\partial \beta_\perp(s)} w_\perp(t) = -\gamma \Theta(t-s) a(s) \quad (116)$$

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

$$\Sigma_{\beta_\star} = (\gamma \mathbf{I} - D_{\beta_\star})^{-1} \kappa_{\beta_\star} (\gamma \mathbf{I} - D_{\beta_\star})^{-1\top}, \quad \Sigma_\perp = (\gamma \mathbf{I} - D_\perp)^{-1} \kappa_\perp (\gamma \mathbf{I} - D_\perp)^{-1\top} \quad (117)$$

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

$$\langle |\beta(t) - \beta_\star|^2 \rangle \sim (\beta(t) - \beta_\star)^2 + \frac{1}{N} \Sigma_{\beta_\star}(t, t) + \frac{(D-1)}{N} \Sigma_{\beta_\perp}(t, t) + \mathcal{O}(N^{-2}) \quad (118)$$

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

### 885 J.3 Connections to Offline Learning in Linear Model

886 **Remark 1** *The finite size variance of generalization error in an online learning setting with linear*  
887 *target function  $y = \beta^\star \cdot \mathbf{x}$  has an identical form as the model described above. In this setting, we*  
888 *sample infinitely many fresh data points  $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I})$  at each step leading to the flow  $\frac{d}{dt} \mathbf{w}_i(t) =$*   
889  *$\gamma a_i(t) \mathbb{E}_{\mathbf{x}} \Delta(\mathbf{x}) \mathbf{x}$  and  $\frac{d}{dt} a_i(t) = \gamma \mathbf{w}_i(t) \cdot \mathbb{E}_{\mathbf{x}} \Delta(\mathbf{x}) \mathbf{x}$ . The order parameter of interest in this setting is*  
890  *$\beta(t) = \frac{1}{\gamma N} \sum_{i=1}^N \mathbf{w}_i(t) a_i(t)$ . The precise correspondence between this setting and the offline setting*  
891 *is summarized in Table I. We note that this argument could be extended to higher degree monomial*  
892 *activations as well, at the cost of tracking higher degree tensors (eg for quadratic activations*  
893  *$\mathbf{M} = \frac{1}{N} \sum_{i=1}^N a_i \mathbf{w}_i \mathbf{w}_i^\top \in \mathbb{R}^{D \times D}$  is sufficient).*

Setting	Order Params.	Target	Off-target Dims.	Loss	Variance	Infinite Quantity
Offline	$\Delta = \mathbf{y} - \mathbf{f}$	$\mathbf{y}$	$P - 1$	Train	$\mathcal{O}(\frac{P}{N})$	$D$
Online	$\beta_\star - \beta$	$\beta_\star$	$D - 1$	Test	$\mathcal{O}(\frac{D}{N})$	$P$

Table 1: Summary of the equivalence between the leading  $1/N$  correction in the offline setting and the online setting for two layer linear networks. In the offline training setting, the order parameters are the errors  $\Delta = \mathbf{y} - \mathbf{f} \in \mathbb{R}^P$  while in the online case they are  $\beta_\star - \beta \in \mathbb{R}^D$ .

894 As in the offline case, in Fig. 3(c) and (d) we see that the variance contribution to test loss  $|\beta - \beta_\star|^2$   
895 increases with input dimension  $D$ . We note that this perturbative effect to the loss dynamics is  
896 reminiscent of the deviations from mean field behavior studied in SGD [43, 44], though this present  
897 work concerns fluctuations driven by initialization variance rather than stochastic sampling of data.  
898 In Fig. 3(e) we show that richer networks have lower variance at fixed  $N$ . Similarly, leading order  
899 theory for richer networks more accurately captures their dynamics as  $D/N$  increases (Fig. 3(f)).

## 900 K Deep Linear Networks

901 For deep linear networks, the fields  $h_\mu^\ell(t), g_\mu^\ell(t)$  are Gaussian and have the following self-consistent  
902 equations

$$\begin{aligned} h_\mu^\ell(t) &= u_\mu^\ell(t) + \gamma \int_0^t ds \sum_\nu [A_{\mu\nu}^{\ell-1}(t, s) + \Delta_\nu(s) H_{\mu\nu}^{\ell-1}(t, s)] g_\nu^\ell(s), \quad u_\mu^\ell(t) \sim \mathcal{GP}(0, \mathbf{H}^{\ell-1}) \\ g_\mu^\ell(t) &= r_\mu^\ell(t) + \gamma \int_0^t ds \sum_\nu [B_{\mu\nu}^\ell(t, s) + \Delta_\nu(s) G_{\mu\nu}^{\ell+1}(t, s)] h_\nu^\ell(s), \quad r_\mu^\ell(t) \sim \mathcal{GP}(0, \mathbf{G}^{\ell+1}). \end{aligned} \quad (119)$$

903 where  $H_{\mu\nu}^\ell(t, s) = \langle h_\mu^\ell(t) h_\nu^\ell(s) \rangle$  and  $G_{\mu\nu}^\ell(t, s) = \langle g_\mu^\ell(t) g_\nu^\ell(s) \rangle$  and  $A_{\mu\nu}^\ell(t, s) = \left\langle \frac{\partial h_\mu^\ell(t)}{\partial r_\nu(s)} \right\rangle$  and

904  $B_{\mu\nu}^\ell(t, s) = \left\langle \frac{\partial h_\mu^\ell(t)}{\partial r_\nu(s)} \right\rangle$  [9]. Therefore, we express the action as a differentiable function of the

order parameters by integrating over the Gaussian field distribution. For concreteness, we vectorize our fields over time and samples  $\mathbf{h}^\ell = \text{Vec}\{h_\mu^\ell(t)\}_{\{\mu \in [P], t \in \mathbb{R}_+\}}$ ,  $\mathbf{g}^\ell = \text{Vec}\{g_\mu^\ell(t)\}_{\{\mu \in [P], t \in \mathbb{R}_+\}}$  we consider the contribution of a single hidden layer.

$$\mathcal{Z}_\ell = \int d\hat{\mathbf{h}}^\ell d\hat{\mathbf{g}}^\ell d\mathbf{h}^\ell d\mathbf{g}^\ell \exp \left( -\frac{1}{2} \hat{\mathbf{h}}^\ell \Sigma_u \hat{\mathbf{h}}^\ell + i \hat{\mathbf{h}}^\ell \cdot (\mathbf{h}^\ell - \mathbf{C}^\ell \mathbf{g}^\ell) - \frac{1}{2} \mathbf{h}^{\ell\top} \hat{\mathbf{H}}^\ell \mathbf{h}^\ell \right) \exp \left( -\frac{1}{2} \hat{\mathbf{g}}^\ell \Sigma_r \hat{\mathbf{g}}^\ell + i \hat{\mathbf{g}}^\ell \cdot (\mathbf{g}^\ell - \mathbf{D}^\ell \mathbf{h}^\ell) - \frac{1}{2} \mathbf{g}^{\ell\top} \hat{\mathbf{G}}^\ell \mathbf{g}^\ell \right)$$

where  $C_{\mu\nu}^\ell(t, s) = \gamma \Theta(t - s) [A_{\mu\nu}^{\ell-1}(t, s) + H_{\mu\nu}^{\ell-1}(t, s) \Delta_\nu(s)]$  and  $D_{\mu\nu}^\ell(t, s) = \gamma \Theta(t - s) [B_{\mu\nu}^\ell(t, s) + G_{\mu\nu}^{\ell+1}(t, s) \Delta_\nu(s)]$ . Performing the joint Gaussian integrals over  $(\mathbf{h}^\ell, \mathbf{g}^\ell, \hat{\mathbf{h}}^\ell, \hat{\mathbf{g}}^\ell)$  we find

$$\ln \mathcal{Z}_\ell = -\frac{1}{2} \ln \det \begin{bmatrix} -\hat{\mathbf{H}}^\ell & 0 & \mathbf{I} & -\mathbf{D}^{\ell\top} \\ 0 & -\hat{\mathbf{G}}^\ell & -\mathbf{C}^{\ell\top} & \mathbf{I} \\ \mathbf{I} & -\mathbf{C}^\ell & \Sigma_u & 0 \\ -\mathbf{D}^\ell & \mathbf{I} & 0 & \Sigma_r \end{bmatrix} \quad (120)$$

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

$$S = \frac{1}{2} \text{Tr} \left[ \hat{\mathbf{H}}^1 \mathbf{H}^1 + \hat{\mathbf{H}}^2 \mathbf{H}^2 + \hat{\mathbf{G}}^1 \mathbf{G}^1 + \hat{\mathbf{G}}^2 \mathbf{G}^2 \right] - \gamma^2 \text{Tr} \mathbf{A} \mathbf{B} - \frac{1}{2} \ln \det \begin{bmatrix} -\hat{\mathbf{H}}^1 & 0 & \mathbf{I} & -\mathbf{D}^{1\top} \\ 0 & -\hat{\mathbf{G}}^1 & -\mathbf{C}^{1\top} & \mathbf{I} \\ \mathbf{I} & -\mathbf{C}^1 & \mathbf{1}\mathbf{1}^\top & 0 \\ -\mathbf{D}^1 & \mathbf{I} & 0 & \mathbf{G}^2 \end{bmatrix} - \frac{1}{2} \ln \det \begin{bmatrix} -\hat{\mathbf{H}}^2 & 0 & \mathbf{I} & -\mathbf{D}^{2\top} \\ 0 & -\hat{\mathbf{G}}^2 & -\mathbf{C}^{2\top} & \mathbf{I} \\ \mathbf{I} & -\mathbf{C}^2 & \mathbf{H}^1 & 0 \\ -\mathbf{D}^2 & \mathbf{I} & 0 & \mathbf{1}\mathbf{1}^\top \end{bmatrix} \quad (121)$$

where  $\mathbf{C}^1 = \gamma \Theta_\Delta$  and  $\mathbf{C}^2 = \gamma \Theta_\Delta \odot \mathbf{H}^1 + \gamma \mathbf{A}$  and  $\mathbf{D}^1 = \gamma \Theta_\Delta \odot \mathbf{G}^2 + \gamma \mathbf{B}$  and  $\mathbf{D}^2 = \gamma \Theta_\Delta$ . 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.

## L 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.

### L.1 Two Layer Linear Equations

In a two layer linear network, the DMFT equations are

$$h(t+1) = h(t) + \eta \gamma \Delta(t) z(t), \quad z(t+1) = z(t) + \eta \gamma \Delta(t) h(t) \\ f(t) = \frac{1}{\gamma} \langle z(t) h(t) \rangle \quad (122)$$

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

$$\begin{aligned} f(t+1) &= f(t) + \eta \langle h(t)^2 + z(t)^2 \rangle \Delta(t) + \eta^2 \gamma \Delta(t)^2 \langle h(t)z(t) \rangle \\ &= f(t) + \eta K(t) \Delta(t) + \eta^2 \gamma^2 \Delta(t)^2 f(t) \end{aligned} \quad (123)$$

$$\begin{aligned} K(t+1) &= K(t) + 4\eta\gamma\Delta(t) \langle h(t)z(t) \rangle + \eta^2 \gamma^2 \Delta(t)^2 \langle h(t)^2 + z(t)^2 \rangle \\ &= K(t) + 4\eta\gamma^2 \Delta(t) f(t) + \eta^2 \gamma^2 \Delta(t)^2 K(t) \end{aligned} \quad (124)$$

930 These equations define the infinite width evolution of  $\Delta(t)$  and  $K(t)$ . Already at this level of analysis,  
 931 we can reason about the evolution of  $K(t)$ . In the small  $\eta$  limit, we could disregard terms of order  
 932  $\mathcal{O}(\eta^2)$  and arrive at the following gradient flow approximation for  $K(t) \sim 2\sqrt{1 + \gamma^2 f(t)^2}$  [9]. This  
 933 evolution will not reach the edge of stability provided that  $\eta < \frac{1}{\sqrt{1 + \gamma^2 y^2}}$ . For large  $\gamma$  and  $y = 1$ , this  
 934 leads to the constraint  $\eta\gamma < 1$ . However, if  $\eta$  exceeds this bound, the gradient flow approximation is  
 935 no longer reasonable and the system reaches an edge of stability effect as shown in Figure 5

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

$$v_+(t+1) = v_+(t) + \eta\gamma\Delta(t)v_+(t), \quad v_-(t+1) = v_-(t) - \eta\gamma\Delta(t)v_-(t). \quad (125)$$

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

$$\begin{aligned} \kappa(t, s) &= 2 \langle h(t)h(s) \rangle^2 + 2 \langle z(t)z(s) \rangle^2 + 2 \langle h(t)z(s) \rangle^2 + 2 \langle z(t)h(s) \rangle^2 \\ &= \langle v_+(t)v_+(s) + v_-(t)v_-(s) \rangle^2 + \langle v_+(t)v_+(s) - v_-(t)v_-(s) \rangle^2. \end{aligned} \quad (126)$$

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

$$\begin{aligned} D(t, s) &= 2 \left\langle v_+(t) \frac{\partial v_+(t)}{\partial \Delta(s)} \right\rangle + 2 \left\langle v_-(t) \frac{\partial v_-(t)}{\partial \Delta(s)} \right\rangle \\ \frac{\partial v_+(t)}{\partial \Delta(s)} &= \gamma \Theta(t-s) v_+(s) + \sum_{t' < t} \Delta(t') \frac{\partial v_+(t')}{\partial \Delta(s)} \\ \frac{\partial v_-(t)}{\partial \Delta(s)} &= -\gamma \Theta(t-s) v_-(s) - \sum_{t' < t} \Delta(t') \frac{\partial v_-(t')}{\partial \Delta(s)} \end{aligned} \quad (127)$$

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

## 943 M Computing Details

944 Experiments for Figures 2, 5 and A.1 were conducted on a Google Colab GPU with JAX. Experiments  
 945 for Figures 4, A.5, 6 were performed on a NVIDIA SMX4-A100-80GB GPU. The total compute  
 946 required for all Figures in the paper took around 4 hours.