
Self-Consistent Dynamical Field Theory of Kernel Evolution in Wide Neural Networks

Blake Bordelon & Cengiz Pehlevan

John Paulson School of Engineering and Applied Sciences, Center for Brain Science
Harvard University
Cambridge MA, 02138

blake_bordelon@g.harvard.edu, cpehlevan@g.harvard.edu

Abstract

We analyze feature learning in infinite-width neural networks trained with gradient flow through a self-consistent dynamical field theory. We construct a collection of deterministic dynamical order parameters which are inner-product kernels for hidden unit activations and gradients in each layer at pairs of time points, providing a reduced description of network activity through training. These kernel order parameters collectively define the hidden layer activation distribution, the evolution of the neural tangent kernel, and consequently output predictions. We show that the field theory derivation recovers the recursive stochastic process of infinite-width feature learning networks obtained from Yang & Hu with Tensor Programs [1]. For deep linear networks, these kernels satisfy a set of algebraic matrix equations. For nonlinear networks, we provide an alternating sampling procedure to self-consistently solve for the kernel order parameters. We provide comparisons of the self-consistent solution to various approximation schemes including the static NTK approximation, gradient independence assumption, and leading order perturbation theory, showing that each of these approximations can break down in regimes where general self-consistent solutions still provide an accurate description. Lastly, we provide experiments in more realistic settings which demonstrate that the loss and kernel dynamics of CNNs at fixed feature learning strength is preserved across different widths on a CIFAR classification task.

1 Introduction

Deep learning has emerged as a successful paradigm for solving challenging machine learning and computational problems across a variety of domains [2, 3]. However, theoretical understanding of the training and generalization of modern deep learning methods lags behind current practice. Ideally, a theory of deep learning would be analytically tractable, efficiently computable, capable of predicting network performance and internal features that the network learns, and interpretable through a reduced description involving desirably initialization-independent quantities.

Several recent theoretical advances have fruitfully considered the idealization of *wide neural networks*, where the number of hidden units in each layer is taken to be large. Under certain parameterization, Bayesian neural networks and gradient descent trained networks converge to gaussian processes (NNGPs) [4-6] and neural tangent kernel (NTK) machines [7-9] in their respective infinite-width limits. These limits provide both analytic tractability as well as detailed training and generalization analysis [10-17]. However, in this limit, with these parameterizations, data representations are fixed and do not adapt to data, termed the *lazy regime* of NN training, to contrast it from the *rich regime* where NNs significantly alter their internal features while fitting the data [18, 19]. The fact that the representation of data is fixed renders these kernel-based theories incapable of explaining feature

learning, an ingredient which is crucial to the success of deep learning in practice [20, 21]. Thus, alternative theories capable of modeling feature learning dynamics are needed.

Recently developed alternative parameterizations such as the mean field [22] and the μP [1] parameterizations allow feature learning in infinite-width NNs trained with gradient descent. Using the Tensor Programs framework, Yang & Hu identified a stochastic process that describes the evolution of preactivation features in infinite-width μP NNs [1]. In this work, we study an equivalent parameterization to μP with self-consistent dynamical mean field theory (DMFT) and recover the stochastic process description of infinite NNs using this alternative technique. In the same large width scaling, we include a scalar parameter γ_0 that allows smooth interpolation between lazy and rich behavior [18]. We provide a new computational procedure to sample this stochastic process and demonstrate its predictive power for wide NNs.

Our novel contributions in this paper are the following:

1. We develop a path integral formulation of gradient flow dynamics in infinite-width networks in the feature learning regime. Our parameterization includes a scalar parameter γ_0 to allow interpolation between rich and lazy regimes and comparison to perturbative methods.
2. Using a stationary action argument, we identify a set of saddle point equations that the kernels satisfy at infinite-width, relating the stochastic processes that define hidden activation evolution to the kernels and vice versa. We show that our saddle point equations recover at $\gamma_0 = 1$, from an alternative method, the same stochastic process obtained previously with Tensor Programs [1].
3. We develop a polynomial-time numerical procedure to solve the saddle point equations for deep networks. In numerical experiments, we demonstrate that solutions to these self-consistency equations are predictive of network training at a variety of feature learning strengths, widths and depths. We provide comparisons of our theory to various approximate methods, such as perturbation theory.

1.1 Related Works

A natural extension to the lazy NTK/NGP limit that allows the study of feature learning is to calculate finite width corrections to the infinite-width limit. Finite width corrections to Bayesian inference in wide networks have been obtained with various perturbative [23-29] and self-consistent techniques [30-33]. In the gradient descent based setting, leading order corrections to the NTK dynamics have been analyzed to study finite width effects [34-36, 27]. These methods give approximate corrections which are accurate provided the strength of feature learning is small. In very rich feature learning regimes, however, the leading order corrections can give incorrect predictions [37, 38].

Another approach to study feature learning is to alter NN parameterization in gradient-based learning to allow significant feature evolution even at infinite-width, the *mean field* limit [22, 39]. Works on mean field NNs have yielded formal loss convergence results [40, 41] and shown equivalences of gradient flow dynamics to a partial differential equation (PDE) [42-44].

Our results are most closely related to a set of recent works which studied infinite-width NNs trained with gradient descent (GD) using the Tensor Programs (TP) framework [1]. We show that our discrete time field theory at unit feature learning strength $\gamma_0 = 1$ recovers the stochastic process which was derived from TP. The stochastic process derived from TP has provided insights into practical issues in NN training such as hyper-parameter search [45]. Computing the exact infinite-width limit of GD has exponential time requirements [1], which we show can be circumvented with an alternating sampling procedure. A projected variant of GD training has provided an infinite-width theory that could be scaled to realistic datasets like CIFAR-10 [46]. Inspired by Chizat and Bach's work on mechanisms of lazy and rich training [18], our theory interpolates between lazy and rich behavior in the mean field limit for varying γ_0 and allows comparison of DMFT to perturbative analysis near small γ_0 . Further, our derivation of a DMFT action allows the possibility of pursuing finite width effects.

Our theory is inspired by self-consistent dynamical mean field theory (DMFT) from statistical physics [47-53]. This framework has been utilized in the theory of random recurrent networks [54-59], tensor PCA [60, 61], phase retrieval [62], and high-dimensional linear classifiers [63-66], but has yet to be developed for deep feature learning. By developing a self-consistent DMFT of deep NNs, we gain insight into how features evolve in the rich regime of network training, while retaining many pleasant analytic properties of the infinite-width limit.

2 Problem Setup and Definitions

Our theory applies to infinite-width networks, both fully-connected and convolutional. For notational ease we will relegate convolutional results to later sections. For input $\mathbf{x}_\mu \in \mathbb{R}^D$, we define the hidden *pre-activation* vectors $\mathbf{h}^\ell \in \mathbb{R}^N$ for layers $\ell \in \{1, \dots, L\}$ as

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

where $\boldsymbol{\theta} = \text{Vec}\{\mathbf{W}^0, \dots, \mathbf{w}^L\}$ are the trainable parameters of the network and ϕ is a twice differentiable activation function. Inspired by previous works on the mechanisms of lazy gradient based training, the parameter γ will control the laziness or richness of the training dynamics [18, 19, 1, 42]. Each of the trainable parameters are initialized as Gaussian random variables with unit variance $W_{ij}^\ell \sim \mathcal{N}(0, 1)$. They evolve under gradient flow $\frac{d}{dt}\boldsymbol{\theta} = -\gamma^2 \nabla_{\boldsymbol{\theta}} \mathcal{L}$. The choice of learning rate γ^2 causes $\frac{d}{dt}\mathcal{L}|_{t=0}$ to be independent of γ . To characterize the evolution of weights, we introduce backpropagation variables $\mathbf{g}_\mu^\ell = \gamma\sqrt{N} \frac{\partial f_\mu}{\partial \mathbf{h}_\mu^\ell} = \dot{\phi}(\mathbf{h}_\mu^\ell) \odot \mathbf{z}_\mu^\ell$, where $\mathbf{z}_\mu^\ell = \frac{1}{\sqrt{N}}\mathbf{W}^{\ell\top} \mathbf{g}_\mu^{\ell+1}$ is the *pre-gradient* signal.

The relevant dynamical objects to characterize feature learning are feature and gradient kernels for each hidden layer $\ell \in \{1, \dots, L\}$, defined as

$$\Phi_{\mu\alpha}^\ell(t, s) = \frac{1}{N} \phi(\mathbf{h}_\mu^\ell(t)) \cdot \phi(\mathbf{h}_\alpha^\ell(s)), \quad G_{\mu\alpha}^\ell(t, s) = \frac{1}{N} \mathbf{g}_\mu^\ell(t) \cdot \mathbf{g}_\alpha^\ell(s). \quad (2)$$

From the kernels $\{\Phi^\ell, G^\ell\}_{\ell=1}^L$, we can compute the *Neural Tangent Kernel* $K_{\mu\alpha}^{NTK}(t, s) = \nabla_{\boldsymbol{\theta}} f_\mu(t) \cdot \nabla_{\boldsymbol{\theta}} f_\alpha(s) = \sum_{\ell=0}^L G_{\mu\alpha}^{\ell+1}(t, s) \Phi_{\mu\alpha}^\ell(t, s)$, [7] and the dynamics of the network function f_μ

$$\frac{d}{dt} f_\mu(t) = \sum_{\alpha=1}^P K_{\mu\alpha}^{NTK}(t, t) \Delta_\alpha(t), \quad \Delta_\mu(t) = -\frac{\partial}{\partial f_\mu} \mathcal{L}|_{f_\mu(t)}, \quad (3)$$

where we define base cases $G_{\mu\alpha}^{L+1}(t, s) = 1$, $\Phi_{\mu\alpha}^0(t, s) = K_{\mu\alpha}^x = \frac{1}{D} \mathbf{x}_\mu \cdot \mathbf{x}_\alpha$. In prior work, Φ^ℓ, G^ℓ were termed *forward* and *backward* kernels and were theoretically computed at initialization and empirically measured through training [67]. Our DMFT will provide exact formulas for these kernels throughout the full dynamics of feature learning. We note that the above formula holds for any data point μ which may or may not be in the set of P training examples. The above expressions demonstrate that knowledge of the temporal trajectory of the NTK on the $t = s$ diagonal gives the temporal trajectory of the network predictions $f_\mu(t)$.

Following prior works on infinite-width networks [22, 1, 40, 19], we study the mean field limit

$$N, \gamma \rightarrow \infty, \quad \gamma_0 = \frac{\gamma}{\sqrt{N}} = \mathcal{O}_N(1) \quad (4)$$

As we demonstrate in the Appendix D and N, this is the only N -scaling which allows feature learning as $N \rightarrow \infty$. The $\gamma_0 = 0$ limit recovers the static NTK limit [7]. We discuss other scalings and parameterizations in Appendix N, relating our work to the μP -parameterization and TP analysis of [1], showing they have identical feature dynamics in the infinite-width limit. We also analyze the effect of different hidden layer widths and initialization variances in the Appendix D.8. We focus on equal widths and NTK parameterization (as in eq. (1)) in the main text to reduce complexity.

3 Self-consistent DMFT

Next, we derive our self-consistent DMFT in a limit where $t, P = \mathcal{O}_N(1)$. Our goal is to build a description of training dynamics purely based on representations, and independent of weights. Studying feature learning at infinite-width enjoys several analytical properties:

- The kernel order parameters Φ^ℓ, G^ℓ concentrate over random initializations but are dynamical, allowing flexible adaptation of features to the task structure.
- In each layer ℓ , each neuron's preactivation h_i^ℓ and pregradient z_i^ℓ become i.i.d. draws from a distribution characterized by a set of order parameters $\{\Phi^\ell, G^\ell, A^\ell, B^\ell\}$.
- The kernels are defined as self-consistent averages (denoted by $\langle \cdot \rangle$) over this distribution of neurons in each layer $\Phi_{\mu\alpha}^\ell(t, s) = \langle \phi(h_\mu^\ell(t)) \phi(h_\alpha^\ell(s)) \rangle$ and $G_{\mu\alpha}^\ell(t, s) = \langle g_\mu^\ell(t) g_\alpha^\ell(s) \rangle$.

The next section derives these facts from a path-integral formulation of gradient flow dynamics.

3.1 Path Integral Construction

Gradient flow after a random initialization of weights defines a high dimensional stochastic process over initializations for variables $\{\mathbf{h}, \mathbf{g}\}$. Therefore, we will utilize DMFT formalism to obtain a reduced description of network activity during training. For a simplified derivation of the DMFT for the two-layer ($L = 1$) case, see [D.2](#). Generally, we separate the contribution on each forward/backward pass between the initial condition and gradient updates to weight matrix \mathbf{W}^ℓ , defining new stochastic variables $\chi^\ell, \xi^\ell \in \mathbb{R}^N$ as

$$\chi_\mu^{\ell+1}(t) = \frac{1}{\sqrt{N}} \mathbf{W}^\ell(0) \phi(\mathbf{h}_\mu^\ell(t)), \quad \xi_\mu^\ell(t) = \frac{1}{\sqrt{N}} \mathbf{W}^\ell(0)^\top \mathbf{g}_\mu^{\ell+1}(t). \quad (5)$$

We let Z represent the moment generating functional (MGF) for these stochastic fields

$$Z[\{\mathbf{j}^\ell, \mathbf{v}^\ell\}] = \left\langle \exp \left(\sum_{\ell, \mu} \int_0^\infty dt [\mathbf{j}_\mu^\ell(t) \cdot \chi_\mu^\ell(t) + \mathbf{v}_\mu^\ell(t) \cdot \xi_\mu^\ell(t)] \right) \right\rangle_{\{\mathbf{W}^0(0), \dots, \mathbf{W}^L(0)\}},$$

which requires, by construction the normalization condition $Z[\{\mathbf{0}, \mathbf{0}\}] = 1$. We enforce our definition of χ, ξ using an integral representation of the delta-function. Thus for each sample $\mu \in [P]$ and each time $t \in \mathbb{R}_+$, we multiply Z by

$$1 = \int_{\mathbb{R}^N} \int_{\mathbb{R}^N} \frac{d\chi_\mu^{\ell+1}(t) d\hat{\chi}_\mu^{\ell+1}(t)}{(2\pi)^N} \exp \left(i \hat{\chi}_\mu^{\ell+1}(t) \cdot \left[\chi_\mu^{\ell+1}(t) - \frac{1}{\sqrt{N}} \mathbf{W}^\ell(0) \phi(\mathbf{h}_\mu^\ell(t)) \right] \right), \quad (6)$$

for χ and the respective expression for ξ . After making such substitutions, we perform integration over initial Gaussian weight matrices to arrive at an integral expression for Z , which we derive in the appendix [D.4](#). We show that Z can be described by set of order-parameters $\{\Phi^\ell, \hat{\Phi}^\ell, G^\ell, \hat{G}^\ell, A^\ell, B^\ell\}$

$$\begin{aligned} Z[\{\mathbf{j}^\ell, \mathbf{v}^\ell\}] &\propto \int \prod_{\ell, \mu, \alpha, t, s} d\Phi_{\mu\alpha}^\ell(t, s) d\hat{\Phi}_{\mu\alpha}^\ell(t, s) dG_{\mu\alpha}^\ell(t, s) d\hat{G}_{\mu\alpha}^\ell(t, s) dA_{\mu\alpha}^\ell(t, s) dB_{\mu\alpha}^\ell(t, s) \quad (7) \\ &\quad \times \exp \left(NS[\{\Phi, \hat{\Phi}, G, \hat{G}, A, B, j, v\}] \right), \\ S &= \sum_{\ell, \mu, \alpha} \int_0^\infty dt \int_0^\infty ds \left[\Phi_{\mu\alpha}^\ell(t, s) \hat{\Phi}_{\mu\alpha}^\ell(t, s) + G_{\mu\alpha}^\ell(t, s) \hat{G}_{\mu\alpha}^\ell(t, s) - A_{\mu\alpha}^\ell(t, s) B_{\mu\alpha}^\ell(t, s) \right] \\ &\quad + \ln \mathcal{Z}[\{\Phi, \hat{\Phi}, G, \hat{G}, A, B, j, v\}], \quad (8) \end{aligned}$$

where S is the DMFT action and \mathcal{Z} is a single-site MGF, which defines the distribution of fields $\{\chi^\ell, \xi^\ell\}$ over the neural population in each layer. The kernels A and B are related to the correlations between feedforward and feedback signals in the network. We provide a detailed formula for \mathcal{Z} in the Appendix [D.4](#) and show that it factorizes over different layers $\mathcal{Z} = \prod_{\ell=1}^L \mathcal{Z}_\ell$.

3.2 Deriving the DMFT Equations from the Path Integral Saddle Point

As $N \rightarrow \infty$, the moment-generating function Z is exponentially dominated by the saddle point of S . The equations that define this saddle point also define our DMFT. We thus identify the kernels that render S locally stationary ($\delta S = 0$). The most important equations are those which define $\{\Phi^\ell, G^\ell\}$

$$\begin{aligned} \frac{\delta S}{\delta \hat{\Phi}_{\mu\alpha}^\ell(t, s)} &= \Phi_{\mu\alpha}^\ell(t, s) + \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta \hat{\Phi}_{\mu\alpha}^\ell(t, s)} = \Phi_{\mu\alpha}^\ell(t, s) - \langle \phi(h_\mu^\ell(t)) \phi(h_\alpha^\ell(s)) \rangle = 0, \\ \frac{\delta S}{\delta \hat{G}_{\mu\alpha}^\ell(t, s)} &= G_{\mu\alpha}^\ell(t, s) + \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta \hat{G}_{\mu\alpha}^\ell(t, s)} = G_{\mu\alpha}^\ell(t, s) - \langle g_\mu^\ell(t) g_\alpha^\ell(s) \rangle = 0, \quad (9) \end{aligned}$$

where $\langle \rangle$ denotes an average over the stochastic process induced by \mathcal{Z} , which is defined below

$$\begin{aligned} \{u_\mu^\ell(t)\}_{\mu \in [P], t \in \mathbb{R}_+} &\sim \mathcal{GP}(0, \Phi^{\ell-1}), \quad \{r_\mu^\ell(t)\}_{\mu \in [P], t \in \mathbb{R}_+} \sim \mathcal{GP}(0, \mathbf{G}^{\ell+1}), \\ h_\mu^\ell(t) &= u_\mu^\ell(t) + \gamma_0 \int_0^t ds \sum_{\alpha=1}^P [A_{\mu\alpha}^{\ell-1}(t, s) + \Delta_\alpha(s) \Phi_{\mu\alpha}^{\ell-1}(t, s)] z_\alpha^\ell(s) \dot{\phi}(h_\alpha^\ell(s)), \\ z_\mu^\ell(t) &= r_\mu^\ell(t) + \gamma_0 \int_0^t ds \sum_{\alpha=1}^P [B_{\mu\alpha}^\ell(t, s) + \Delta_\alpha(s) G_{\mu\alpha}^{\ell+1}(t, s)] \phi(h_\alpha^\ell(s)), \quad (10) \end{aligned}$$

where we define base cases $\Phi_{\mu\alpha}^0(t, s) = K_{\mu\alpha}^x$ and $G_{\mu\alpha}^{L+1}(t, s) = 1$, $A^0 = B^L = 0$. We see that the fields $\{h^\ell, z^\ell\}$, which represent the single site preactivations and pre-gradients, are implicit functionals of the mean-zero Gaussian processes $\{u^\ell, r^\ell\}$ which have covariances $\langle u_\mu^\ell(t)u_\alpha^\ell(s) \rangle = \Phi_{\mu\alpha}^{\ell-1}(t, s)$ and $\langle r_\mu^\ell(t)r_\alpha^\ell(s) \rangle = G_{\mu\alpha}^{\ell+1}(t, s)$. The other saddle point equations give $A_{\mu\alpha}^\ell(t, s) = \gamma_0^{-1} \left\langle \frac{\delta\phi(h_\mu^\ell(t))}{\delta r_\alpha^\ell(s)} \right\rangle$, $B_{\mu\alpha}^\ell(t, s) = \gamma_0^{-1} \left\langle \frac{\delta g_\mu^{\ell+1}(t)}{\delta u_\alpha^{\ell+1}(s)} \right\rangle$ which arise due to coupling between the feedforward and feedback signals. We note that, in the lazy limit $\gamma_0 \rightarrow 0$, the fields approach Gaussian processes $h^\ell \rightarrow u^\ell$, $z^\ell \rightarrow r^\ell$. Lastly, the final saddle point equations $\frac{\delta S}{\delta \Phi^\ell} = 0$, $\frac{\delta S}{\delta G^\ell} = 0$ imply that $\hat{\Phi}^\ell = \hat{G}^\ell = 0$. The full set of equations that define the DMFT are given in [D.7](#).

This theory is easily extended to more general architectures such as networks with varying widths by layer (App. [D.8](#)), trainable bias parameter (App. [H](#)), multiple (but $\mathcal{O}_N(1)$) output channels (App. [I](#)), convolutional architectures (App. [G](#)), networks trained with weight decay (App. [J](#)), Langevin sampling (App. [K](#)) and momentum (App. [L](#)), discrete time training (App. [M](#)). In Appendix [N](#), we discuss parameterizations which give equivalent feature and predictor dynamics and show our derived stochastic process is equivalent to the μP scheme of Yang & Hu [II](#).

4 Solving the Self-Consistent DMFT

The saddle point equations obtained from the field theory discussed in the previous section must be solved self-consistently. By this we mean that, given knowledge of the kernels, we can characterize the distribution of $\{h^\ell, z^\ell\}$, and given the distribution of $\{h^\ell, z^\ell\}$, we can compute the kernels [\[68, 64\]](#). In the Appendix [B](#), we provide Algorithm [I](#), a numerical procedure based on this idea to efficiently solve for the kernels with an alternating Monte-Carlo strategy. The output of the algorithm are the dynamical kernels $\Phi_{\mu\alpha}^\ell(t, s)$, $G_{\mu\alpha}^\ell(t, s)$, $A_{\mu\alpha}^\ell(t, s)$, $B_{\mu\alpha}^\ell(t, s)$, from which any network observable can be computed as we discuss in Appendix [D](#). We provide an example of the solution to the saddle point equations compared to training a finite NN in Figure [I](#). We plot Φ^ℓ, G^ℓ at the end of training and the sample-trace of these kernels through time. Additionally, we compare the kernels of finite width N network to the DMFT predicted kernels using a cosine-similarity alignment metric $A(\Phi^{DMFT}, \Phi^{NN}) = \frac{\text{Tr} \Phi^{DMFT} \Phi^{NN}}{|\Phi^{DMFT}| |\Phi^{NN}|}$. Additional examples are in Appendix Figures [6](#) and Figure [7](#).

4.1 Deep Linear Networks: Closed Form Self-Consistent Equations

Deep linear networks ($\phi(h) = h$) are of theoretical interest since they are simpler to analyze than nonlinear networks but preserve non-trivial training dynamics and feature learning [\[69-73, 25, 32, 23\]](#). In a deep linear network, we can simplify our saddle point equations to algebraic formulas that close in terms of the kernels $H_{\mu\alpha}^\ell(t, s) = \langle h_\mu^\ell(t)h_\alpha^\ell(s) \rangle$, $G^\ell(t, s) = \langle g^\ell(t)g^\ell(s) \rangle$ [\[II\]](#). This is a significant simplification since it allows solution of the saddle point equations without a sampling procedure.

To describe the result, we first introduce a vectorization notation $\mathbf{h}^\ell = \text{Vec}\{h_\mu^\ell(t)\}_{\mu \in [P], t \in \mathbb{R}_+}$. Likewise we convert kernels $\mathbf{H}^\ell = \text{Mat}\{H_{\mu\alpha}^\ell(t, s)\}_{\mu, \alpha \in [P], t, s \in \mathbb{R}_+}$ into matrices. The inner product under this vectorization is defined as $\mathbf{a} \cdot \mathbf{b} = \int_0^\infty dt \sum_{\mu=1}^P a_\mu(t)b_\mu(t)$. In a practical computational implementation, the theory would be evaluated on a grid of T time points with discrete time gradient descent, so these kernels $\mathbf{H}^\ell \in \mathbb{R}^{PT \times PT}$ would indeed be matrices of the appropriate size. The fields $\mathbf{h}^\ell, \mathbf{g}^\ell$ are linear functionals of independent Gaussian processes $\mathbf{u}^\ell, \mathbf{r}^\ell$, giving $(\mathbf{I} - \gamma_0^2 \mathbf{C}^\ell \mathbf{D}^\ell) \mathbf{h}^\ell = \mathbf{u}^\ell + \gamma_0 \mathbf{C}^\ell \mathbf{r}^\ell$, $(\mathbf{I} - \gamma_0^2 \mathbf{D}^\ell \mathbf{C}^\ell) \mathbf{g}^\ell = \mathbf{r}^\ell + \gamma_0 \mathbf{D}^\ell \mathbf{u}^\ell$. The matrices \mathbf{C}^ℓ and \mathbf{D}^ℓ are causal integral operators which depend on $\{A^{\ell-1}, \mathbf{H}^{\ell-1}\}$ and $\{B^\ell, \mathbf{G}^{\ell+1}\}$ respectively which we define in Appendix [F](#). The saddle point equations which define the kernels are

$$\begin{aligned} \mathbf{H}^\ell &= \langle \mathbf{h}^\ell \mathbf{h}^{\ell\top} \rangle = (\mathbf{I} - \gamma_0^2 \mathbf{C}^\ell \mathbf{D}^\ell)^{-1} [\mathbf{H}^{\ell-1} + \gamma_0^2 \mathbf{C}^\ell \mathbf{G}^{\ell+1} \mathbf{C}^{\ell\top}] [(\mathbf{I} - \gamma_0^2 \mathbf{C}^\ell \mathbf{D}^\ell)^{-1}]^\top \\ \mathbf{G}^\ell &= \langle \mathbf{g}^\ell \mathbf{g}^{\ell\top} \rangle = (\mathbf{I} - \gamma_0^2 \mathbf{D}^\ell \mathbf{C}^\ell)^{-1} [\mathbf{G}^{\ell+1} + \gamma_0^2 \mathbf{D}^\ell \mathbf{H}^{\ell-1} \mathbf{D}^{\ell\top}] [(\mathbf{I} - \gamma_0^2 \mathbf{D}^\ell \mathbf{C}^\ell)^{-1}]^\top. \end{aligned} \quad (11)$$

Examples of the predictions obtained by solving these systems of equations are provided in Figure [2](#). We see that these DMFT equations describe kernel evolution for networks of a variety of depths and that the change in each layer's kernel increases with the depth of the network.

Unlike many prior results [\[69-72\]](#), our DMFT does not require any restrictions on the structure of the input data but hold for any \mathbf{K}^x, \mathbf{y} . However, for whitened data $\mathbf{K}^x = \mathbf{I}$ we show in

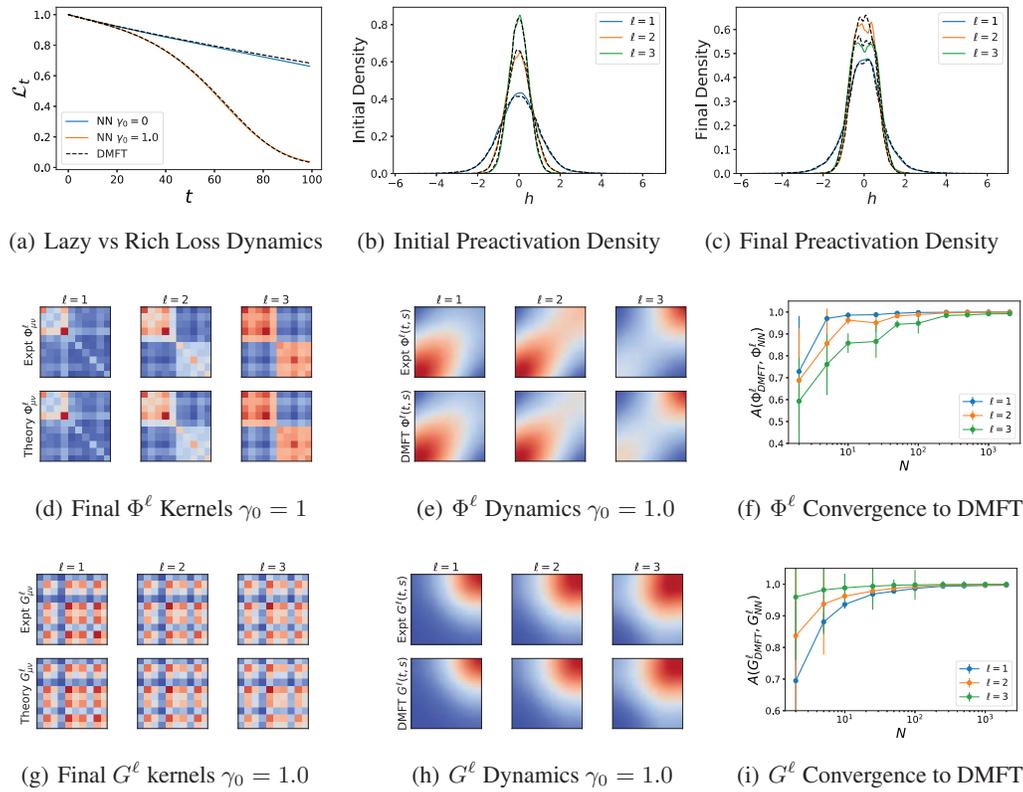


Figure 1: Neural network feature learning dynamics is captured by self-consistent dynamical mean field theory (DMFT). (a) Training loss curves on a subsample of $P = 10$ CIFAR-10 training points in a depth 4 ($L = 3$, $N = 2500$) tanh network ($\phi(h) = \tanh(h)$) trained with MSE. Increasing γ_0 accelerates training. (b)-(c) The distribution of preactivations at the beginning and end of training matches predictions of the DMFT. (d) The final Φ^ℓ (at $t = 100$) kernel order parameters match the finite width network. (e) The temporal dynamics of the sample-traced kernels $\sum_\mu \Phi_{\mu\mu}^\ell(t, s)$ matches experiment and reveals rich dynamics across layers. (f) The alignment $A(\Phi_{DMFT}^\ell, \Phi_{NN}^\ell)$, defined as cosine similarity, of the kernel $\Phi_{\mu\alpha}^\ell(t, s)$ predicted by theory (DMFT) and width N networks for different N but fixed $\gamma_0 = \gamma/\sqrt{N}$. Errorbars show standard deviation computed over 10 repeats. Around $N \sim 500$ DMFT begins to show near perfect agreement with the NN. (g)-(i) The same plots but for the gradient kernel G^ℓ . Whereas finite width effects for Φ^ℓ are larger at later layers ℓ since variance accumulates on the forward pass, fluctuations in G^ℓ are large in early layers.

Appendix [F.1.1](#), [F.2](#) that our DMFT learning curves interpolate between NTK dynamics and the sigmoidal trajectories of prior works [\[69, 70\]](#) as γ_0 is increased. For example, in the two layer ($L = 1$) linear network with $K^x = \mathbf{I}$, the dynamics of the error norm $\Delta(t) = \|\Delta(t)\|$ takes the form $\frac{\partial}{\partial t} \Delta(t) = -2\sqrt{1 + \gamma_0^2(y - \Delta(t))^2} \Delta(t)$ where $y = \|\mathbf{y}\|$. These dynamics give the linear convergence rate of the NTK if $\gamma_0 \rightarrow 0$ but approaches logistic dynamics of [\[70\]](#) as $\gamma_0 \rightarrow \infty$. Further, $\mathbf{H}(t) = \langle \mathbf{h}^1(t) \mathbf{h}^1(t)^\top \rangle \in \mathbb{R}^{P \times P}$ only grows in the $\mathbf{y}\mathbf{y}^\top$ direction with $H_y(t) = \frac{1}{y^2} \mathbf{y}^\top \mathbf{H}(t) \mathbf{y} = \sqrt{1 + \gamma_0^2(y - \Delta(t))^2}$. At the end of training $\mathbf{H}(t) \rightarrow \mathbf{I} + \frac{1}{y^2} [\sqrt{1 + \gamma_0^2 y^2} - 1] \mathbf{y}\mathbf{y}^\top$, recovering the rank one spike which was recently obtained in the small initialization limit [\[74\]](#). We show this one dimensional system in Figure [8](#).

4.2 Feature Learning with L2 Regularization

As we show in Appendix [D](#), the DMFT can be extended to networks trained with weight decay $\frac{d\theta}{dt} = -\gamma^2 \nabla_\theta \mathcal{L} - \lambda \theta$. If neural network is homogenous in its parameters so that $f(c\theta) = c^\kappa f(\theta)$ (examples include networks with linear, ReLU, quadratic activations), then the final network predictor

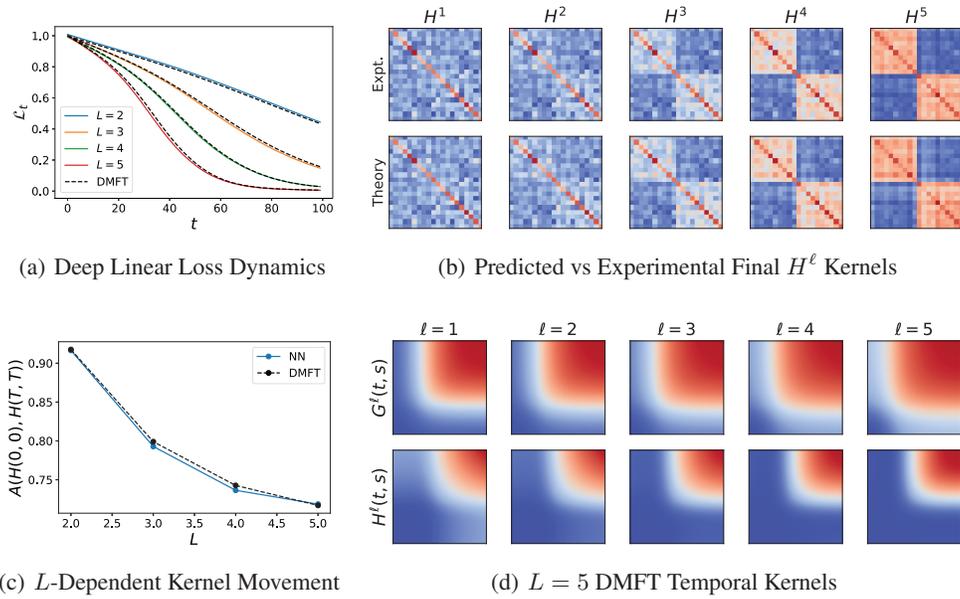


Figure 2: Deep linear network with the full DMFT. (a) The train loss for NNs of varying L . (b) For a $L = 5$, $N = 1000$ NN, the kernels H^ℓ at the end of training compared to DMFT theory on $P = 20$ datapoints. (c) The average displacement of feature kernels for different depth networks at same γ_0 value. For equal values of γ_0 , deeper networks exhibit larger changes to their features, manifested in lower alignment with their initial $t = 0$ kernels H . (d) The solution to the temporal components of the $G^\ell(t, s)$ and $\sum_\mu H_{\mu\mu}^\ell(t, s)$ kernels obtained from the self-consistent equations.

is a kernel regressor with the final NTK $\lim_{t \rightarrow \infty} f(\mathbf{x}, t) = \mathbf{k}(\mathbf{x})^\top [\mathbf{K} + \lambda\kappa\mathbf{I}]^{-1} \mathbf{y}$ where $K(\mathbf{x}, \mathbf{x}')$ is the *final*-NTK, $[\mathbf{k}(\mathbf{x})]_\mu = K(\mathbf{x}, \mathbf{x}_\mu)$ and $[\mathbf{K}]_{\mu\alpha} = K(\mathbf{x}_\mu, \mathbf{x}_\alpha)$. We note that the effective regularization $\lambda\kappa$ increases with depth L . In NTK parameterization, weight decay in infinite width homogenous networks gives a trivial fixed point $K(\mathbf{x}, \mathbf{x}') \rightarrow 0$ and consequently a zero predictor $f \rightarrow 0$ [75]. However, as we show in Figure 3, increasing feature learning γ_0 can prevent convergence to the trivial fixed point, allowing a non-zero fixed point for K, f even at infinite width. The kernel and function dynamics can be predicted with DMFT. The fixed point is a nontrivial function of the hyperparameters $\lambda, \kappa, L, \gamma_0$.

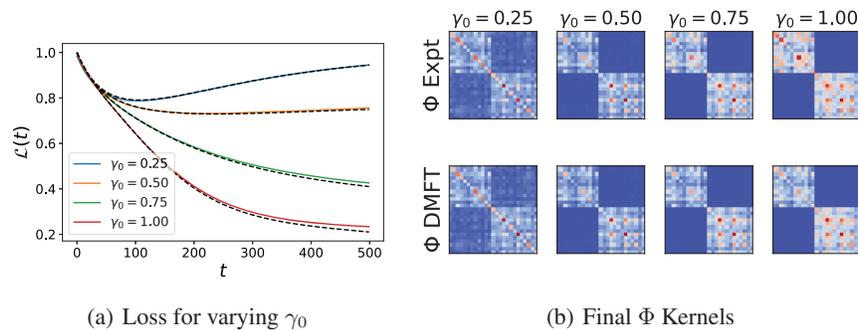


Figure 3: Width $N = 1000$ ReLU networks trained with L2 regularization have nontrivial fixed point in DMFT limit ($\gamma_0 > 0$). (a) Training loss dynamics for a $L = 1$ ReLU network with $\lambda = 1$. In $\gamma_0 \rightarrow 0$ limit the fixed point is trivial $f = K = 0$. The final loss is a decreasing function of γ_0 . (b) The final kernel is more aligned with target with increasing γ_0 . Networks with homogenous activations enjoy a representer theorem at infinite-width as we show in Appendix J.

5 Approximation Schemes

We now compare our exact DMFT with approximations of prior works, providing an explanation of when these approximations give accurate predictions and when they break down.

5.1 Gradient Independence Ansatz

We can study the accuracy of the ansatz $\mathbf{A}^\ell = \mathbf{B}^\ell = 0$, which is equivalent to treating the weight matrices $\mathbf{W}^\ell(0)$ and $\mathbf{W}^{\ell(0)\top}$ which appear in forward and backward passes respectively as independent Gaussian matrices. This assumption was utilized in prior works on signal propagation in deep networks in the lazy regime [76–80]. A consequence of this approximation is the Gaussianity and statistical independence of χ^ℓ and ξ^ℓ (conditional on $\{\Phi^\ell, \mathbf{G}^\ell\}$) in each layer as we show in Appendix O. This ansatz works very well near $\gamma_0 \approx 0$ (the static kernel regime) since $\frac{dh}{dr}, \frac{dz}{du} \sim \mathcal{O}(\gamma_0)$ or around initialization $t \approx 0$ but begins to fail at larger values of γ_0, t (Figure 4).

5.2 Perturbation theory in γ_0 at infinite-width

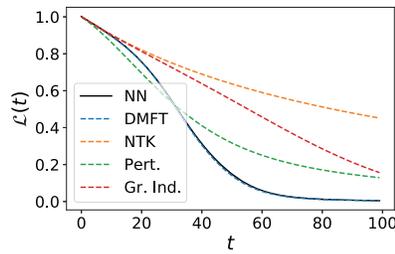
In the $\gamma_0 \rightarrow 0$ limit, we recover static kernels, giving linear dynamics identical to the NTK limit [7]. Corrections to this lazy limit can be extracted at small but finite γ_0 . This is conceptually similar to recent works which consider perturbation series for the NTK in powers of $1/N$ [35, 27, 28] (though not identical, see Appendix P.7 for finite N effects). We expand all observables $q(\gamma_0)$ in a power series in γ_0 , giving $q(\gamma_0) = q^{(0)} + \gamma_0 q^{(1)} + \gamma_0^2 q^{(2)} + \dots$ and compute corrections up to $\mathcal{O}(\gamma_0^2)$. We show that the $\mathcal{O}(\gamma_0)$ and $\mathcal{O}(\gamma_0^3)$ corrections to kernels vanish, giving leading order expansions of the form $\Phi = \Phi^0 + \gamma_0^2 \Phi^2 + \mathcal{O}(\gamma_0^4)$ and $\mathbf{G} = \mathbf{G}^0 + \gamma_0^2 \mathbf{G}^2 + \mathcal{O}(\gamma_0^4)$ (see Appendix P.2). Further, we show that the NTK has relative change at leading order which scales linearly with depth $|\Delta K^{NTK}|/|K^{NTK,0}| \sim \mathcal{O}_{\gamma_0, L}(L\gamma_0^2) = \mathcal{O}_{N, \gamma, L}(\frac{\gamma^2 L}{N})$, which is consistent with finite width effective field theory at $\gamma = \mathcal{O}_N(1)$ [26–28] (Appendix P.6). Further, at the leading order correction, all temporal dependencies are controlled by $P(P+1)$ functions $v_\alpha(t) = \int_0^t ds \Delta_\alpha^0(s)$ and $v_{\alpha\beta}(t) = \int_0^t ds \Delta_\alpha^0(s) \int_0^s ds' \Delta_\beta^0(s')$, which is consistent with those derived for finite width NNs using a truncation of the Neural Tangent Hierarchy [34, 35, 27]. To lighten notation, we focus our main text comparison of our non-perturbative DMFT to perturbation theory in the deep linear case. Full perturbation theory is in Appendix P.2.

Using the timescales derived in the previous section, we find that the leading order correction to the kernels in infinite-width deep linear network have the form

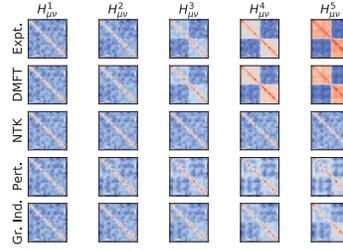
$$K_{\mu\nu}^{NTK}(t, s) = (L+1)K_{\mu\nu}^x + \gamma_0^2 \frac{L(L+1)}{2} K_{\mu\nu}^x \sum_{\alpha\beta} K_{\alpha\beta}^x [v_{\alpha\beta}(t) + v_{\beta\alpha}(s) + v_\alpha(t)v_\beta(s)] \\ + \gamma_0^2 \frac{L(L+1)}{2} \left[\sum_{\alpha\beta} K_{\mu\alpha}^x K_{\nu\beta}^x [v_{\alpha\beta}(t) + v_{\beta\alpha}(s)] + \sum_{\alpha\beta} K_{\mu\alpha}^x K_{\nu\beta}^x v_\alpha(t)v_\beta(s) \right] + \mathcal{O}(\gamma_0^4). \quad (12)$$

We see that the relative change in the NTK $|\mathbf{K}^{NTK} - \mathbf{K}^{NTK}(0)|/|\mathbf{K}^{NTK}(0)| \sim \mathcal{O}(\gamma_0^2 L) = \mathcal{O}(\gamma^2 L/N)$, so that large depth L networks exhibit more significant kernel evolution, which agrees with other perturbative studies [35, 27, 25] as well as the non-perturbative results in Figure 2. However at large γ_0 and large L , this theory begins to break down as we show in Figure 4.

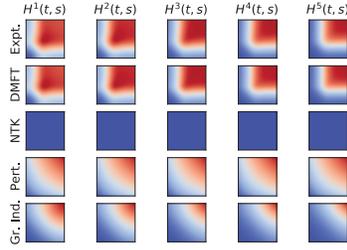
The DMFT formalism can also be used to extract leading corrections to observables at large but finite width N as we explore in P.7. When deviating from infinite width, the kernels are no longer deterministic over network initializations. The key observation is that the DMFT action S defines a Gibbs measure over the space of kernel order parameters $\mathbf{k} = \text{Vec}\{\Phi^\ell, \mathbf{G}^\ell, \mathbf{A}^\ell, \mathbf{B}^\ell\}$ with probability density $\frac{1}{Z} \exp(NS[\mathbf{k}])$ where Z is a normalization constant. Near infinite width, any observable average $\langle O(\mathbf{k}) \rangle = \frac{1}{Z} \int d\mathbf{k} \exp(NS[\mathbf{k}]) O(\mathbf{k})$ is dominated by order parameters within a $\frac{1}{\sqrt{N}}$ neighborhood of \mathbf{k}^* . As a consequence, a perturbative series for $\langle O(\mathbf{k}) \rangle$ can be obtained from simple averages over Gaussian fluctuations in the kernels $\mathbf{k} \sim \mathcal{N}(\mathbf{k}^*, -\frac{1}{N}[\nabla^2 S[\mathbf{k}^*]]^{-1})$ [29]. The components for $\nabla^2 S[\mathbf{k}^*]$ include four point correlations of fields computed over the DMFT distribution.



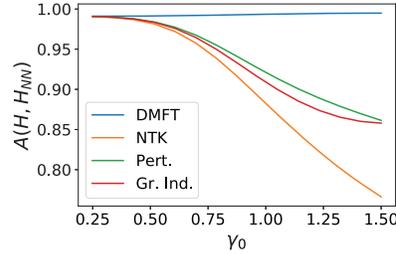
(a) Loss dynamics



(b) Final H^ℓ Kernels $\gamma_0 = 1.5$



(c) H^ℓ Kernel Dynamics $\gamma_0 = 1.5$



(d) Theory H^ℓ vs NN with $N = 1000$

Figure 4: Comparison of DMFT to various approximation schemes in a $L = 5$ hidden layer, width $N = 1000$ linear network with $\gamma_0 = 1.0$ and $P = 100$. (a) The loss for the various approximations do not track the true trajectory induced by gradient descent in the large γ_0 regime. (b)-(c) The feature kernels $H_{\mu\alpha}^\ell(t, s)$ across each of the $L = 5$ hidden layers for each of the theories is compared to a width 1000 neural network. Again, we plot the sample-traced dynamics $\sum_{\mu\mu} H_{\mu\mu}^\ell(t, s)$. (d) The alignment of H^ℓ compared to the finite NN $A(H^\ell, H_{NN}^\ell)$ averaged across $\ell \in \{1, \dots, 5\}$ for varying γ . The predictions of all of these theories coincide in the $\gamma_0 = 0$ limit but begin to deviate in the feature learning regime. Only the non-perturbative DMFT is accurate over a wide range of γ_0 .

6 Feature Learning Dynamics is Preserved at Fixed γ_0

Our DMFT suggests that for networks sufficiently wide for their kernels to concentrate, the dynamics of loss and kernels should be invariant under the rescaling $N \rightarrow RN, \gamma \rightarrow \gamma/\sqrt{R}$, which keeps γ_0 fixed. To evaluate how well this idea holds in a realistic deep learning problem, we trained CNNs of varying channel counts N on two-class CIFAR classification [81]. We tracked the dynamics of the loss and the last layer Φ^L kernel. The results are provided in Figure 5. We see that dynamics are largely independent of rescaling as predicted. Further, as expected, larger γ_0 leads to larger changes in kernel norm and faster alignment to the target function y , as was also found in [82]. Consequently, the higher γ_0 networks train more rapidly. The trend is consistent for width $N = 250$ and $N = 500$. More details about the experiment can be found in Appendix C.2.

7 Discussion

We provided a unifying DMFT derivation of feature dynamics in infinite networks trained with gradient based optimization. Our theory interpolates between lazy infinite-width behavior of a static NTK in $\gamma_0 \rightarrow 0$ and rich feature learning. At $\gamma_0 = 1$, our DMFT construction agrees with the stochastic process derived previously with the Tensor Programs framework [1]. Our saddle point equations give self-consistency conditions which relate the stochastic fields to the kernels. These equations are exactly solvable in deep linear networks and can be efficiently solved with a numerical method in the nonlinear case. Comparisons with other approximation schemes show that DMFT can be accurate at a much wider range of γ_0 . We believe our framework could be a useful perspective for future theoretical analyses of feature learning and generalization in wide networks.

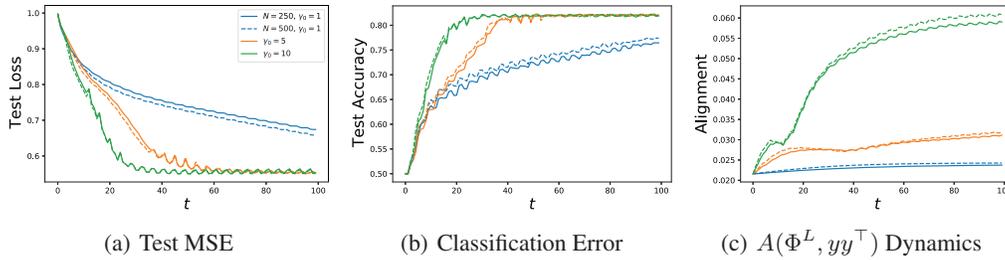


Figure 5: The dynamics of a depth 5 ($L = 4$ hidden) CNNs trained on first two classes of CIFAR (boat vs plane) exhibit consistency for different channel counts $N \in \{250, 500\}$ for fixed $\gamma_0 = \gamma/\sqrt{N}$. (a) We plot the test loss (MSE) and (b) test classification error. Networks with higher γ_0 train more rapidly. Time is measured in every 100 update steps. (c) The dynamics of the last layer feature kernel Φ^L , shown as alignment to the target function. As predicted by the DMFT, higher γ_0 corresponds to more active kernel evolution, evidenced by larger change in the alignment.

Though our DMFT is quite general in regards to the data and architecture, the technique is not entirely rigorous and relies on heuristic physics techniques. Our theory holds in the $T, P = \mathcal{O}_N(1)$ and may break down otherwise; other asymptotic regimes (such as $P/N, T/\log(N) = \mathcal{O}_N(1)$, etc) may exhibit phenomena relevant to deep learning practice [32, 83]. The computational requirements of our method, while smaller than the exponential time complexity for exact solution [11], are still significant for large PT . In Table 1, we compare the time taken for various theories to compute the feature kernels throughout T steps of gradient descent. For a width N network, computation of each forward pass on all P data points takes $\mathcal{O}(PN^2)$ computations. The static NTK requires computation of $\mathcal{O}(P^2)$ entries in the kernel which do not need to be recomputed. However, the DMFT requires matrix multiplications on $PT \times PT$ matrices giving a $\mathcal{O}(P^3T^3)$ time scaling. Future work could aim to improve the computational overhead of the algorithm, by considering data averaged theories [64] or one pass SGD [11]. Alternative projected versions of gradient descent have also enabled much better computational scaling in evaluation of the theoretical predictions [46], allowing evaluation on full CIFAR-10.

Requirements	Width- N NN	Static NTK	Perturbative	Full DMFT
Memory for Kernels	$\mathcal{O}(N^2)$	$\mathcal{O}(P^2)$	$\mathcal{O}(P^4T)$	$\mathcal{O}(P^2T^2)$
Time for Kernels	$\mathcal{O}(PN^2T)$	$\mathcal{O}(P^2)$	$\mathcal{O}(P^4T)$	$\mathcal{O}(P^3T^3)$
Time for Final Outputs	$\mathcal{O}(PN^2T)$	$\mathcal{O}(P^3)$	$\mathcal{O}(P^4)$	$\mathcal{O}(P^3T^3)$

Table 1: Computational requirements to compute kernel dynamics and trained network predictions on P points in a depth N neural network on a grid of T time points trained with P data points for various theories. DMFT is faster and less memory intensive than a width N network only if $N \gg PT$. It is more computationally efficient to compute full DMFT kernels than leading order perturbation theory when $T \ll \sqrt{P}$. The expensive scaling with both samples and time are the cost of a full-batch non-perturbative theory of gradient based feature learning dynamics.

Acknowledgments and Disclosure of Funding

This work was supported by NSF grant DMS-2134157 and an award from the Harvard Data Science Initiative Competitive Research Fund. BB acknowledges additional support from the NSF-Simons Center for Mathematical and Statistical Analysis of Biology at Harvard (award #1764269) and the Harvard Q-Bio Initiative.

BB thanks Jacob Zavatore-Veth, Alex Atanasov, Abdulkadir Canatar, and Ben Ruben for comments on this manuscript as well as Greg Yang, Boris Hanin, Yasaman Bahri, and Jascha Sohl-Dickstein for useful discussions.

References

- [1] Greg Yang and Edward J Hu. Tensor programs iv: Feature learning in infinite-width neural networks. In *International Conference on Machine Learning*, pages 11727–11737. PMLR, 2021.
- [2] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- [3] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- [4] Radford M Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media, 2012.
- [5] Jaehoon Lee, Jascha Sohl-dickstein, Jeffrey Pennington, Roman Novak, Sam Schoenholz, and Yasaman Bahri. Deep neural networks as gaussian processes. In *International Conference on Learning Representations*, 2018.
- [6] Alexander G. de G. Matthews, Jiri Hron, Mark Rowland, Richard E. Turner, and Zoubin Ghahramani. Gaussian process behaviour in wide deep neural networks. In *International Conference on Learning Representations*, 2018.
- [7] Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31, pages 8571–8580. Curran Associates, Inc., 2018.
- [8] Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. *Advances in neural information processing systems*, 32, 2019.
- [9] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. *Advances in Neural Information Processing Systems*, 32, 2019.
- [10] Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *International conference on machine learning*, pages 1675–1685. PMLR, 2019.
- [11] B. Bordelon, A. Canatar, and C. Pehlevan. Spectrum dependent learning curves in kernel regression and wide neural networks. *International Conference of Machine Learning*, 2020.
- [12] Abdulkadir Canatar, Blake Bordelon, and Cengiz Pehlevan. Spectral bias and task-model alignment explain generalization in kernel regression and infinitely wide neural networks. *Nature communications*, 12(1):1–12, 2021.
- [13] Omry Cohen, Or Malka, and Zohar Ringel. Learning curves for overparametrized deep neural networks: A field theory perspective. *Physical Review Research*, 3(2):023034, 2021.
- [14] Arthur Jacot, Berfin Simsek, Francesco Spadaro, Clément Hongler, and Franck Gabriel. Kernel alignment risk estimator: Risk prediction from training data. *Advances in Neural Information Processing Systems*, 33:15568–15578, 2020.
- [15] Bruno Loureiro, Cedric Gerbelot, Hugo Cui, Sebastian Goldt, Florent Krzakala, Marc Mezard, and Lenka Zdeborova. Learning curves of generic features maps for realistic datasets with a teacher-student model. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.
- [16] James B Simon, Madeline Dickens, and Michael R DeWeese. Neural tangent kernel eigenvalues accurately predict generalization. *arXiv preprint arXiv:2110.03922*, 2021.
- [17] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In *International Conference on Machine Learning*, pages 242–252. PMLR, 2019.

- [18] Lenaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. *Advances in Neural Information Processing Systems*, 32, 2019.
- [19] Mario Geiger, Stefano Spigler, Arthur Jacot, and Matthieu Wyart. Disentangling feature and lazy training in deep neural networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2020(11):113301, 2020.
- [20] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [22] Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671, 2018.
- [23] Laurence Aitchison. Why bigger is not always better: on finite and infinite neural networks. In *International Conference on Machine Learning*, pages 156–164. PMLR, 2020.
- [24] Sho Yaida. Non-gaussian processes and neural networks at finite widths. In *Mathematical and Scientific Machine Learning*, pages 165–192. PMLR, 2020.
- [25] Jacob Zavatone-Veth, Abdulkadir Canatar, Ben Ruben, and Cengiz Pehlevan. Asymptotics of representation learning in finite bayesian neural networks. *Advances in Neural Information Processing Systems*, 34, 2021.
- [26] Gadi Naveh, Oded Ben David, Haim Sompolinsky, and Zohar Ringel. Predicting the outputs of finite deep neural networks trained with noisy gradients. *Physical Review E*, 104(6):064301, 2021.
- [27] Daniel A Roberts, Sho Yaida, and Boris Hanin. The principles of deep learning theory. *arXiv preprint arXiv:2106.10165*, 2021.
- [28] Boris Hanin. Correlation functions in random fully connected neural networks at finite width. *arXiv preprint arXiv:2204.01058*, 2022.
- [29] Kai Segadlo, Bastian Epping, Alexander van Meegen, David Dahmen, Michael Krämer, and Moritz Helias. Unified field theory for deep and recurrent neural networks, 2021.
- [30] Gadi Naveh and Zohar Ringel. A self consistent theory of gaussian processes captures feature learning effects in finite cnns. *Advances in Neural Information Processing Systems*, 34, 2021.
- [31] Inbar Seroussi and Zohar Ringel. Separation of scales and a thermodynamic description of feature learning in some cnns. *arXiv preprint arXiv:2112.15383*, 2021.
- [32] Qianyi Li and Haim Sompolinsky. Statistical mechanics of deep linear neural networks: The backpropagating kernel renormalization. *Physical Review X*, 11(3):031059, 2021.
- [33] Jacob A Zavatone-Veth and Cengiz Pehlevan. Depth induces scale-averaging in overparameterized linear bayesian neural networks. *55th Asilomar Conference on Signals, Systems, and Computers*, 2021.
- [34] Jiaoyang Huang and Horng-Tzer Yau. Dynamics of deep neural networks and neural tangent hierarchy. In *International conference on machine learning*, pages 4542–4551. PMLR, 2020.
- [35] Ethan Dyer and Guy Gur-Ari. Asymptotics of wide networks from feynman diagrams. *arXiv preprint arXiv:1909.11304*, 2019.
- [36] Anders Andreassen and Ethan Dyer. Asymptotics of wide convolutional neural networks. *arXiv preprint arXiv:2008.08675*, 2020.

- [37] Jacob A Zavatone-Veth, William L Tong, and Cengiz Pehlevan. Contrasting random and learned features in deep bayesian linear regression. *arXiv preprint arXiv:2203.00573*, 2022.
- [38] Aitor Lewkowycz, Yasaman Bahri, Ethan Dyer, Jascha Sohl-Dickstein, and Guy Gur-Ari. The large learning rate phase of deep learning: the catapult mechanism. *arXiv preprint arXiv:2003.02218*, 2020.
- [39] Dyego Araújo, Roberto I Oliveira, and Daniel Yukimura. A mean-field limit for certain deep neural networks. *arXiv preprint arXiv:1906.00193*, 2019.
- [40] Lenaic Chizat and Francis Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. *Advances in neural information processing systems*, 31, 2018.
- [41] Grant M Rotskoff and Eric Vanden-Eijnden. Trainability and accuracy of neural networks: An interacting particle system approach. *arXiv preprint arXiv:1805.00915*, 2018.
- [42] Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit. In *Conference on Learning Theory*, pages 2388–2464. PMLR, 2019.
- [43] Phan-Minh Nguyen. Mean field limit of the learning dynamics of multilayer neural networks. *arXiv preprint arXiv:1902.02880*, 2019.
- [44] Cong Fang, Jason Lee, Pengkun Yang, and Tong Zhang. Modeling from features: a mean-field framework for over-parameterized deep neural networks. In *Conference on learning theory*, pages 1887–1936. PMLR, 2021.
- [45] Greg Yang, Edward Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ryder, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tuning large neural networks via zero-shot hyperparameter transfer. *Advances in Neural Information Processing Systems*, 34, 2021.
- [46] Greg Yang, Michael Santacroce, and Edward J Hu. Efficient computation of deep nonlinear infinite-width neural networks that learn features. In *International Conference on Learning Representations*, 2022.
- [47] Paul Cecil Martin, ED Siggia, and HA Rose. Statistical dynamics of classical systems. *Physical Review A*, 8(1):423, 1973.
- [48] C De Dominicis. Dynamics as a substitute for replicas in systems with quenched random impurities. *Physical Review B*, 18(9):4913, 1978.
- [49] Haim Sompolinsky and Annette Zippelius. Dynamic theory of the spin-glass phase. *Physical Review Letters*, 47(5):359, 1981.
- [50] Haim Sompolinsky and Annette Zippelius. Relaxational dynamics of the edwards-anderson model and the mean-field theory of spin-glasses. *Physical Review B*, 25(11):6860, 1982.
- [51] G Ben Arous and Alice Guionnet. Large deviations for langevin spin glass dynamics. *Probability Theory and Related Fields*, 102(4):455–509, 1995.
- [52] G Ben Arous and Alice Guionnet. Symmetric langevin spin glass dynamics. *The Annals of Probability*, 25(3):1367–1422, 1997.
- [53] Gérard Ben Arous, Amir Dembo, and Alice Guionnet. Cugliandolo-kurchan equations for dynamics of spin-glasses. *Probability theory and related fields*, 136(4):619–660, 2006.
- [54] A Crisanti and H Sompolinsky. Path integral approach to random neural networks. *Physical Review E*, 98(6):062120, 2018.
- [55] Haim Sompolinsky, Andrea Crisanti, and Hans-Jurgen Sommers. Chaos in random neural networks. *Physical review letters*, 61(3):259, 1988.

- [56] Moritz Helias and David Dahmen. *Statistical Field Theory for Neural Networks*. Springer International Publishing, 2020.
- [57] Lutz Molgedey, J Schuchhardt, and Heinz G Schuster. Suppressing chaos in neural networks by noise. *Physical review letters*, 69(26):3717, 1992.
- [58] M Samuelides and Bruno Cessac. Random recurrent neural networks dynamics. *The European Physical Journal Special Topics*, 142(1):89–122, 2007.
- [59] Kanaka Rajan, LF Abbott, and Haim Sompolinsky. Stimulus-dependent suppression of chaos in recurrent neural networks. *Physical review e*, 82(1):011903, 2010.
- [60] Stefano Sarao Mannelli, Florent Krzakala, Pierfrancesco Urbani, and Lenka Zdeborova. Passed & spurious: Descent algorithms and local minima in spiked matrix-tensor models. In *international conference on machine learning*, pages 4333–4342. PMLR, 2019.
- [61] Stefano Sarao Mannelli, Giulio Biroli, Chiara Cammarota, Florent Krzakala, Pierfrancesco Urbani, and Lenka Zdeborová. Marvels and pitfalls of the langevin algorithm in noisy high-dimensional inference. *Physical Review X*, 10(1):011057, 2020.
- [62] Francesca Mignacco, Pierfrancesco Urbani, and Lenka Zdeborová. Stochasticity helps to navigate rough landscapes: comparing gradient-descent-based algorithms in the phase retrieval problem. *Machine Learning: Science and Technology*, 2(3):035029, 2021.
- [63] Elisabeth Agoritsas, Giulio Biroli, Pierfrancesco Urbani, and Francesco Zamponi. Out-of-equilibrium dynamical mean-field equations for the perceptron model. *Journal of Physics A: Mathematical and Theoretical*, 51(8):085002, 2018.
- [64] Francesca Mignacco, Florent Krzakala, Pierfrancesco Urbani, and Lenka Zdeborová. Dynamical mean-field theory for stochastic gradient descent in gaussian mixture classification. *Advances in Neural Information Processing Systems*, 33:9540–9550, 2020.
- [65] Michael Celentano, Chen Cheng, and Andrea Montanari. The high-dimensional asymptotics of first order methods with random data. *arXiv preprint arXiv:2112.07572*, 2021.
- [66] Francesca Mignacco and Pierfrancesco Urbani. The effective noise of stochastic gradient descent. *arXiv preprint arXiv:2112.10852*, 2021.
- [67] Yizhang Lou, Chris E Mingard, and Soufiane Hayou. Feature learning and signal propagation in deep neural networks. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 14248–14282. PMLR, 17–23 Jul 2022.
- [68] Alessandro Manacorda, Grégory Schehr, and Francesco Zamponi. Numerical solution of the dynamical mean field theory of infinite-dimensional equilibrium liquids. *The Journal of chemical physics*, 152(16):164506, 2020.
- [69] Kenji Fukumizu. Dynamics of batch learning in multilayer neural networks. In *International Conference on Artificial Neural Networks*, pages 189–194. Springer, 1998.
- [70] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- [71] Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient descent for deep linear neural networks. In *International Conference on Learning Representations*, 2019.
- [72] Madhu S Advani, Andrew M Saxe, and Haim Sompolinsky. High-dimensional dynamics of generalization error in neural networks. *Neural Networks*, 132:428–446, 2020.
- [73] Arthur Jacot, François Ged, Franck Gabriel, Berfin Şimşek, and Clément Hongler. Deep linear networks dynamics: Low-rank biases induced by initialization scale and l2 regularization. *arXiv preprint arXiv:2106.15933*, 2021.

- [74] Alexander Atanasov, Blake Bordelon, and Cengiz Pehlevan. Neural networks as kernel learners: The silent alignment effect. In *International Conference on Learning Representations*, 2022.
- [75] Aitor Lewkowycz and Guy Gur-Ari. On the training dynamics of deep networks with l_2 regularization. *Advances in Neural Information Processing Systems*, 33:4790–4799, 2020.
- [76] Ben Poole, Subhaneil Lahiri, Maithra Raghu, Jascha Sohl-Dickstein, and Surya Ganguli. Exponential expressivity in deep neural networks through transient chaos. *Advances in neural information processing systems*, 29, 2016.
- [77] Samuel S Schoenholz, Justin Gilmer, Surya Ganguli, and Jascha Sohl-Dickstein. Deep information propagation. *International Conference of Learning Representations*, 2017.
- [78] Greg Yang and Samuel Schoenholz. Mean field residual networks: On the edge of chaos. *Advances in neural information processing systems*, 30, 2017.
- [79] Greg Yang. Scaling limits of wide neural networks with weight sharing: Gaussian process behavior, gradient independence, and neural tangent kernel derivation. *arXiv preprint arXiv:1902.04760*, 2019.
- [80] Greg Yang and Etai Littwin. Tensor programs iib: Architectural universality of neural tangent kernel training dynamics. In *International Conference on Machine Learning*, pages 11762–11772. PMLR, 2021.
- [81] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [82] Haozhe Shan and Blake Bordelon. A theory of neural tangent kernel alignment and its influence on training, 2021.
- [83] Stéphane d’Ascoli, Maria Refinetti, and Giulio Biroli. Optimal learning rate schedules in high-dimensional non-convex optimization problems, 2022.
- [84] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018.
- [85] Juha Honkonen. Ito and stratonovich calculus in stochastic field theory. *arXiv preprint arXiv:1102.1581*, 2011.
- [86] Crispin W Gardiner et al. *Handbook of stochastic methods*, volume 3. springer Berlin, 1985.
- [87] Carl M Bender and Steven Orszag. *Advanced mathematical methods for scientists and engineers I: Asymptotic methods and perturbation theory*, volume 1. Springer Science & Business Media, 1999.
- [88] John Hubbard. Calculation of partition functions. *Physical Review Letters*, 3(2):77, 1959.
- [89] Charles Stein. A bound for the error in the normal approximation to the distribution of a sum of dependent random variables. In *Proceedings of the sixth Berkeley symposium on mathematical statistics and probability, volume 2: Probability theory*, volume 6, pages 583–603. University of California Press, 1972.
- [90] Roman Novak, Lechao Xiao, Jaehoon Lee, Yasaman Bahri, Greg Yang, Jiri Hron, Daniel A Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein. Bayesian deep convolutional networks with many channels are gaussian processes. *arXiv preprint arXiv:1810.05148*, 2018.
- [91] Greg Yang. Wide feedforward or recurrent neural networks of any architecture are gaussian processes. *Advances in Neural Information Processing Systems*, 32, 2019.
- [92] Adam X. Yang, Maxime Robeyns, Edward Milsom, Nandi Schoots, and Laurence Aitchison. A theory of representation learning in deep neural networks gives a deep generalisation of kernel methods, 2021.

- [93] Yurii E Nesterov. A method for solving the convex programming problem with convergence rate $O(1/k^2)$. In *Dokl. akad. nauk Sssr*, volume 269, pages 543–547, 1983.
- [94] Yurii Nesterov and Boris T Polyak. Cubic regularization of newton method and its global performance. *Mathematical Programming*, 108(1):177–205, 2006.
- [95] Gabriel Goh. Why momentum really works. *Distill*, 2017.
- [96] Michael Muehlebach and Michael I Jordan. Optimization with momentum: Dynamical, control-theoretic, and symplectic perspectives. *Journal of Machine Learning Research*, 22(73):1–50, 2021.
- [97] Mehran Kardar. *Statistical physics of fields*. Cambridge University Press, 2007.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes] As described in the abstract and introduction, we provide a dynamical field theory of deep networks based on kernel evolution.
 - (b) Did you describe the limitations of your work? [Yes] We have an explicit limitations as the last paragraph of the paper in Section 7.
 - (c) Did you discuss any potential negative societal impacts of your work? [N/A] This work is theoretical and is very unlikely to present negative social impacts.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes] We describe that our theory holds for NN architectures in the infinite-width $N \rightarrow \infty$ limit.
 - (b) Did you include complete proofs of all theoretical results? [Yes] All claims made in the main text are supported by derivations in the Appendix.
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] Code to reproduce experimental results is provided in the supplementary material.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] We provide details of all experiments in C.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] We provided errorbars in the alignment scores of DMFT as a function of width N in Figure 1. All other runs were over a single wide network, where performance is predicted to concentrate over initialization.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] We mention our GPU usage in C.2
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes] We cited the creators of Jax, Neural Tangents, and CIFAR-10.
 - (b) Did you mention the license of the assets? [N/A] These are all open source provided they are appropriately credited in academic research.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]

- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]