



The End of Training

The job of a scientist is to listen carefully to nature, not to tell nature how to behave.

Freeman Dyson, explaining Richard Feynman's approach [69].

In this chapter, we'll finally finish our leading-order effective-theory analysis of finite-width networks and solve their training dynamics under gradient descent. In contrast to the infinite-width limit, for which the solution is independent of the training algorithm, the dynamics of such deep networks have a rich phenomenology that captures the different ways in which useful features may develop over the course of training. The solution to these training dynamics gives a first-principles description of the ensemble of fully-trained finite-width networks, realizing a main goal of the book.

Unfortunately, our job will be disrupted by two facts of nature: *(i)* in order to have a consistent description of training dynamics at order $1/n$, we'll need to incorporate two additional objects that arise in the Taylor expansion of the update to the network output to third order, the update to the NTK to second order, and the update to the dNTK to first order; and *(ii)* due to a lack of smoothness, we won't be able to describe the dynamics of ReLU networks or networks consisting of any of the other nonlinear activation functions from the scale-invariant universality class.

As for the first point, while the analysis of representation learning in the context of the quadratic model was illuminating, we've already telegraphed that it was insufficient to capture the particular details of finite-width networks. In particular, to leading order in $1/n$, there are two more NTK differentials, which we'll refer to as *ddNTKs*. Although it's straightforward, working out the stochastic forward equations, recursions, and effective theory for these ddNTKs is somewhat tedious and no longer has any pedagogical value. As such, we won't provide the details of our derivations – you've already seen these sets of manipulations three times before in §4–§5, §8–§9, and §11.2–§11.3, for the preactivations, NTK, and dNTK, respectively – instead we'll simply state the results, leaving the details for you as a kind of post-training test evaluation; after all, this is the end of your training as well.

As for the second point, throughout the book we've had to use special methods in order to work out exceptional explanations for any nonsmooth activation function, such as the ReLU. In our minds, this extra work was justified by the ReLU's current privileged status as one of the most popular activation functions in practice. However, we have finally run out of tricks and will have to give up – for a reason that is simple to explain: our Taylor expansion in the global learning rate η will break down when applied to the dynamics of networks built with nonsmooth activation functions. Instead, we'll have to follow the direction of the community and begin thinking again about smoothed versions of the ReLU – though only the ones that permit a type of criticality – such as the GELU and the SWISH.

With both those disruptions to our work heard, in § $\infty.1$ we'll present all the relevant results for the ddNTKs – we'll define them, we'll give their tensor decomposition, and we'll explain their scaling with width and depth – while hiding all the irrelevant details at the back of the chapter in § $\infty.3$. If you've been paying attention, you'll not be shocked to hear that – when properly normalized – the ddNTKs scale as the effective theory cutoff: ℓ/n . This scaling indicates that we need to consider the joint statistics of the preactivation–NTK–dNTK–ddNTKs in order to understand the leading-order finite-width dynamics of deep MLPs. Importantly, these ddNTKs endow the dNTK with its own dynamics; from the parameter-space perspective of §11.4.1, this means that the *meta feature functions* of the model will now evolve.

With those results stated, in § $\infty.2$ we'll return to our regularly scheduled pedagogy and, at long last, solve the training dynamics at finite width. After an initial false start following our infinite-width giant leap, first in § $\infty.2.1$ we'll learn how to take a small step following an adjusted giant leap, giving us our first finite-width solution. Then in § $\infty.2.2$, we'll analyze many many steps of vanilla gradient descent, giving us our second finite-width solution. The nonlinear dynamics at finite width ultimately lead to a dependence of the fully-trained solution on the training algorithm, and so the solutions derived in these two subsections actually exhibit meaningful differences.

In particular, the function approximation of a fully-trained finite-width network can be decomposed into a universal part, independent of the optimization details, and a set of *algorithm projectors*, whose functional form encodes the entire dependence of the solution on the training algorithm. These projectors provide a dual sample-space perspective on the learning algorithm, analogous to the relationship between the model parameters and the different kernels.

Accordingly, in § $\infty.2.3$ we'll discuss how these projectors impact the solution, letting us understand the inductive bias of the *training dynamics* separately from the inductive bias of the *network architecture*. We'll also further analyze the predictions made by such fully-trained networks, considering the growing tradeoff between increased representation learning and increased instantiation-to-instantiation fluctuations with network depth.

While this is the final chapter of the main text, in a small epilogue following this chapter, Epilogue ε , we'll explore how to define model complexity for overparameterized networks from our effective theory's macroscopic perspective. Then in two appendices, we'll further touch on some topics that are outside the scope of our main line of inquiry. In Appendix A, we'll introduce the framework of information theory, which

will give us the tools we need in order to estimate the optimal aspect ratio that separates *effectively-deep* networks from *overly-deep* networks. In Appendix B, we'll apply our effective theory approach to learn about residual networks and see how they can be used to extend the range of effectively-deep networks to greater and greater depths.

∞.1 Two More Differentials

Who ordered that?

I. I. Rabi, quipping about the $O(1/n)$ ddNTKs.

One last time, let's expand the ℓ -th-layer preactivations after a parameter update, this time recording terms up to *third order*:

$$\begin{aligned} \vec{dz}_{i;\delta}^{(\ell)} &\equiv z_{i;\delta}^{(\ell)}(t=1) - z_{i;\delta}^{(\ell)}(t=0) \\ &= \sum_{\ell_1=1}^{\ell} \sum_{\mu} \frac{dz_{i;\delta}^{(\ell)}}{d\theta_{\mu}^{(\ell_1)}} d\theta_{\mu}^{(\ell_1)} + \frac{1}{2} \sum_{\ell_1, \ell_2=1}^{\ell} \sum_{\mu_1, \mu_2} \frac{d^2 z_{i;\delta}^{(\ell)}}{d\theta_{\mu_1}^{(\ell_1)} d\theta_{\mu_2}^{(\ell_2)}} d\theta_{\mu_1}^{(\ell_1)} d\theta_{\mu_2}^{(\ell_2)} \\ &\quad + \frac{1}{6} \sum_{\ell_1, \ell_2, \ell_3=1}^{\ell} \sum_{\mu_1, \mu_2, \mu_3} \frac{d^3 z_{i;\delta}^{(\ell)}}{d\theta_{\mu_1}^{(\ell_1)} d\theta_{\mu_2}^{(\ell_2)} d\theta_{\mu_3}^{(\ell_3)}} d\theta_{\mu_1}^{(\ell_1)} d\theta_{\mu_2}^{(\ell_2)} d\theta_{\mu_3}^{(\ell_3)} + \dots \end{aligned} \quad (\infty.1)$$

For gradient descent, also recall that after use of the chain rule (11.3), the change in the ℓ_a -th-layer parameters of any particular network is given by (11.4),

$$d\theta_{\mu}^{(\ell_a)} = -\eta \sum_{\nu} \lambda_{\mu\nu}^{(\ell_a)} \left(\sum_{j,k,\tilde{\alpha}} \frac{\partial \mathcal{L}_{\mathcal{A}}}{\partial z_{k;\tilde{\alpha}}^{(L)}} \frac{dz_{k;\tilde{\alpha}}^{(L)}}{dz_{j;\tilde{\alpha}}^{(\ell)}} \frac{dz_{j;\tilde{\alpha}}^{(\ell)}}{d\theta_{\nu}^{(\ell_a)}} \right) = -\eta \sum_{\nu,j,\tilde{\alpha}} \lambda_{\mu\nu}^{(\ell_a)} \epsilon_{j;\tilde{\alpha}}^{(\ell)} \frac{dz_{j;\tilde{\alpha}}^{(\ell)}}{d\theta_{\nu}^{(\ell_a)}}, \quad (\infty.2)$$

where we've used our convention from §11.1 of explicitly specifying which layer each parameter comes from. Please also recall from there that the learning-rate tensor $\lambda_{\mu\nu}^{(\ell)}$ only connects the parameters within a given layer ℓ . In the above expression, ℓ is an intermediate layer such that $\ell_a \leq \ell$, and we also used our ℓ -th-layer error factor (11.5):

$$\epsilon_{j;\tilde{\alpha}}^{(\ell)} \equiv \sum_{k=1}^{n_L} \frac{\partial \mathcal{L}_{\mathcal{A}}}{\partial z_{k;\tilde{\alpha}}^{(L)}} \frac{dz_{k;\tilde{\alpha}}^{(L)}}{dz_{j;\tilde{\alpha}}^{(\ell)}} = \frac{d\mathcal{L}_{\mathcal{A}}}{dz_{j;\tilde{\alpha}}^{(\ell)}}. \quad (\infty.3)$$

After substituting the parameter update (∞.2) back into the preactivation update (∞.1), you should be able to write it in the form

$$\begin{aligned} dz_{i;\delta}^{(\ell)} &= -\eta \sum_{j,\tilde{\alpha}} \widehat{H}_{ij;\delta\tilde{\alpha}}^{(\ell)} \epsilon_{j;\tilde{\alpha}}^{(\ell)} + \frac{\eta^2}{2} \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2} \widehat{dH}_{ij_1j_2;\delta\tilde{\alpha}_1\tilde{\alpha}_2}^{(\ell)} \epsilon_{j_1;\tilde{\alpha}_1}^{(\ell)} \epsilon_{j_2;\tilde{\alpha}_2}^{(\ell)} \\ &\quad - \frac{\eta^3}{6} \sum_{j_1,j_2,j_3,\tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3} \widehat{dd_I H}_{ij_1j_2j_3;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3}^{(\ell)} \epsilon_{j_1;\tilde{\alpha}_1}^{(\ell)} \epsilon_{j_2;\tilde{\alpha}_2}^{(\ell)} \epsilon_{j_3;\tilde{\alpha}_3}^{(\ell)} \\ &\quad + O(\eta^4), \end{aligned} \quad (\infty.4)$$

where the first two terms we found in the last chapter (11.9), and the cubic term is new, with the *first* of the **ddNTKs** defined as

$$\widehat{\text{dd}_I H}_{i_0 i_1 i_2 i_3; \delta_0 \delta_1 \delta_2 \delta_3}^{(\ell)} \equiv \sum_{\ell_1, \ell_2, \ell_3=1}^{\ell} \sum_{\substack{\mu_1, \nu_1, \\ \mu_2, \nu_2, \\ \mu_3, \nu_3}} \lambda_{\mu_1 \nu_1}^{(\ell_1)} \lambda_{\mu_2 \nu_2}^{(\ell_2)} \lambda_{\mu_3 \nu_3}^{(\ell_3)} \frac{d^3 z_{i_0; \delta_0}^{(\ell)}}{d\theta_{\mu_1}^{(\ell_1)} d\theta_{\mu_2}^{(\ell_2)} d\theta_{\mu_3}^{(\ell_3)}} \frac{dz_{i_1; \delta_1}^{(\ell)}}{d\theta_{\nu_1}^{(\ell_1)}} \frac{dz_{i_2; \delta_2}^{(\ell)}}{d\theta_{\nu_2}^{(\ell_2)}} \frac{dz_{i_3; \delta_3}^{(\ell)}}{d\theta_{\nu_3}^{(\ell_3)}}. \quad (\infty.5)$$

As always, the hat on the ddNTK indicates that it's stochastic, depending on the particular realization of the model parameters at initialization. Also, similar to the dNTK, this ddNTK is totally symmetric in its second, third, and fourth paired sets of indices $(i_1, \delta_1) \leftrightarrow (i_2, \delta_2) \leftrightarrow (i_3, \delta_3)$, while the first neural-sample index (i_0, δ_0) is distinguished from the other three.

By expanding to order η^3 , the update, (∞.4), is now *cubic* in error factors. Expanding to this order is necessary because the first ddNTK has statistics at initialization that are $O(1/n)$ and so needs to be included in our analysis. However, any higher-order terms in the update are subleading, so we may replace $O(\eta^4) = O(1/n^2)$ in this expression.

Just as we had to expand the update to the NTK to order η when we expanded the update to the preactivations to order η^2 , we will now have to expand the update to the NTK to order η^2 for the dynamics with our cubic update (∞.4) to be consistent:

$$\begin{aligned} dH_{i_1 i_2; \delta_1 \delta_2}^{(\ell)} &\equiv H_{i_1 i_2; \delta_1 \delta_2}^{(\ell)}(t=1) - H_{i_1 i_2; \delta_1 \delta_2}^{(\ell)}(t=0) \\ &= \sum_{\ell_1=1}^{\ell} \sum_{\mu_1} \frac{dH_{i_1 i_2; \delta_1 \delta_2}^{(\ell)}}{d\theta_{\mu_1}^{(\ell_1)}} d\theta_{\mu_1}^{(\ell_1)} + \sum_{\ell_1, \ell_2=1}^{\ell} \sum_{\mu_1, \mu_2} \frac{d^2 H_{i_1 i_2; \delta_1 \delta_2}^{(\ell)}}{d\theta_{\mu_1}^{(\ell_1)} d\theta_{\mu_2}^{(\ell_2)}} d\theta_{\mu_1}^{(\ell_1)} d\theta_{\mu_2}^{(\ell_2)} + \dots \\ &= -\eta \sum_{j, \tilde{\alpha}} \left(\widehat{dH}_{i_1 i_2 j; \delta_1 \delta_2 \tilde{\alpha}}^{(\ell)} + \widehat{dH}_{i_2 i_1 j; \delta_2 \delta_1 \tilde{\alpha}}^{(\ell)} \right) \epsilon_{j; \tilde{\alpha}}^{(\ell)} \\ &\quad + \frac{\eta^2}{2} \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2} \left[\widehat{\text{dd}_I H}_{i_1 i_2 j_1 j_2; \delta_1 \delta_2 \tilde{\alpha}_1 \tilde{\alpha}_2}^{(\ell)} + \widehat{\text{dd}_I H}_{i_2 i_1 j_1 j_2; \delta_2 \delta_1 \tilde{\alpha}_1 \tilde{\alpha}_2}^{(\ell)} \right] \epsilon_{j_1; \tilde{\alpha}_1}^{(\ell)} \epsilon_{j_2; \tilde{\alpha}_2}^{(\ell)} \\ &\quad + \eta^2 \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2} \widehat{\text{dd}_{II} H}_{i_1 i_2 j_1 j_2; \delta_1 \delta_2 \tilde{\alpha}_1 \tilde{\alpha}_2}^{(\ell)} \epsilon_{j_1; \tilde{\alpha}_1}^{(\ell)} \epsilon_{j_2; \tilde{\alpha}_2}^{(\ell)} + O(\eta^3). \end{aligned} \quad (\infty.6)$$

To go to the final equality, we substituted in our NTK definition (11.7) and our parameter update (∞.2), computed the derivatives, and then collected the terms. To do so, we identified the *second* of the **ddNTKs**, defined as

$$\widehat{\text{dd}_{II} H}_{i_1 i_2 i_3 i_4; \delta_1 \delta_2 \delta_3 \delta_4}^{(\ell)} \equiv \sum_{\ell_1, \ell_2, \ell_3=1}^{\ell} \sum_{\substack{\mu_1, \nu_1, \\ \mu_2, \nu_2, \\ \mu_3, \nu_3}} \lambda_{\mu_1 \nu_1}^{(\ell_1)} \lambda_{\mu_2 \nu_2}^{(\ell_2)} \lambda_{\mu_3 \nu_3}^{(\ell_3)} \frac{d^2 z_{i_1; \delta_1}^{(\ell)}}{d\theta_{\mu_1}^{(\ell_1)} d\theta_{\mu_3}^{(\ell_3)}} \frac{d^2 z_{i_2; \delta_2}^{(\ell)}}{d\theta_{\mu_2}^{(\ell_2)} d\theta_{\nu_3}^{(\ell_3)}} \frac{dz_{i_3; \delta_3}^{(\ell)}}{d\theta_{\nu_1}^{(\ell_1)}} \frac{dz_{i_4; \delta_4}^{(\ell)}}{d\theta_{\nu_2}^{(\ell_2)}}. \quad (\infty.7)$$

The hat on this ddNTK indicates that it's also stochastic at initialization, and we will soon detail that it also has $O(1/n)$ statistics at leading order. Finally, $\widehat{\text{dd}_{II} H}_{i_1 i_2 i_3 i_4; \delta_1 \delta_2 \delta_3 \delta_4}$

has a more constrained symmetry, only symmetric under a joint swap of the paired set of indices as $(i_1, \delta_1) \leftrightarrow (i_2, \delta_2)$ and $(i_3, \delta_4) \leftrightarrow (i_4, \delta_4)$. However, this means that we can also swap indices as

$$\sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2} \widehat{\text{dd}}_{\text{II}} H_{i_1 i_2 j_1 j_2; \delta_1 \delta_2 \tilde{\alpha}_1 \tilde{\alpha}_2}^{(\ell)} \epsilon_{j_1; \tilde{\alpha}_1}^{(\ell)} \epsilon_{j_2; \tilde{\alpha}_2}^{(\ell)} = \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2} \widehat{\text{dd}}_{\text{II}} H_{i_2 i_1 j_1 j_2; \delta_2 \delta_1 \tilde{\alpha}_1 \tilde{\alpha}_2}^{(\ell)} \epsilon_{j_1; \tilde{\alpha}_1}^{(\ell)} \epsilon_{j_2; \tilde{\alpha}_2}^{(\ell)}, \quad (\infty.8)$$

which we used to simplify the final line of the NTK update (∞.6).

Overall, this NTK update is now *quadratic* in error factors, making its dynamics coupled and nonlinear. Again, expanding to this order is necessary because both ddNTKs have statistics at initialization that are $O(1/n)$, and so both need to be included in our analysis. However, any higher-order terms in the NTK update are subleading, so in (∞.6) we may replace $O(\eta^3) = O(1/n^2)$.

Finally, consider the leading-order update to the dNTK:

$$\begin{aligned} \ddot{d}H_{i_0 i_1 i_2; \delta_0 \delta_1 \delta_2}^{(\ell)} &\equiv dH_{i_0 i_1 i_2; \delta_0 \delta_1 \delta_2}^{(\ell)}(t=1) - dH_{i_0 i_1 i_2; \delta_0 \delta_1 \delta_2}^{(\ell)}(t=0) \\ &= \sum_{\ell_1=1}^{\ell} \sum_{\mu_1} \frac{d dH_{i_0 i_1 i_2; \delta_0 \delta_1 \delta_2}^{(\ell)}}{d\theta_{\mu_1}^{(\ell_1)}} \dot{\theta}_{\mu_1}^{(\ell_1)} + \dots \\ &= -\eta \sum_{j, \tilde{\alpha}} \left(\widehat{\text{dd}}_{\text{I}} H_{i_0 i_1 i_2 j; \delta_0 \delta_1 \delta_2 \tilde{\alpha}}^{(\ell)} + \widehat{\text{dd}}_{\text{II}} H_{i_0 i_1 i_2 j; \delta_0 \delta_1 \delta_2 \tilde{\alpha}}^{(\ell)} + \widehat{\text{dd}}_{\text{II}} H_{i_0 i_2 i_1 j; \delta_0 \delta_2 \delta_1 \tilde{\alpha}}^{(\ell)} \right) \epsilon_{j; \tilde{\alpha}}^{(\ell)} \\ &\quad + O(\eta^2). \end{aligned} \quad (\infty.9)$$

To go to the final equality, we substituted our dNTK definition (11.8) and parameter update (∞.2), took the derivative, and then collected all the terms using our ddNTK definitions, (∞.5) and (∞.7). The dNTK update is *linear* in error factors, and its dynamics will be the simplest. Finally, the higher-order terms in the dNTK update are subleading, so in (∞.9) we may replace $O(\eta^2) = O(1/n^2)$. We also would like to apologize for the $\ddot{d}dH$ notation, and we promise that we won't have to use it again.

The updates (∞.4), (∞.6), and (∞.9) comprise the complete set of finite-width updates at $O(1/n)$. Thus, to proceed further in our analysis, we'll need to work out the leading-order statistics of the ddNTKs.

ddNTK Statistics

To find the distribution of fully-trained finite-width networks, we need the joint distribution of the network output, the NTK, the dNTK, and the ddNTKs:

$$p\left(z^{(L)}, \hat{H}^{(L)}, \widehat{dH}^{(L)}, \widehat{\text{dd}}_{\text{I}} H^{(L)}, \widehat{\text{dd}}_{\text{II}} H^{(L)} \middle| \mathcal{D}\right). \quad (\infty.10)$$

Rather than working through the details here, we'll do it in our own private notebooks. You already have all the tools you need: you can follow §4, §8, and §11.2 for examples of how to use RG flow to work out the recursions; and you can follow §5, §9, and §11.3

for examples of how to work out the details of the effective theory at initialization after tuning to criticality. The full details of this distribution (∞.10) can be found at the end of the chapter in §∞.3. Here, we'll highlight the results that you need in order to understand our dynamical computations in the following section.

After working out the stochastic forward equations for both ddNTKs – feel free to flip forward to (∞.174) and (∞.175) if you're curious about them – we'll need to find recursions for its statistics. For these tensors, the leading-order statistics come from their means; their cross correlations and their variances are all subleading. When evaluating the mean of each of the ddNTKs, it will become convenient to decompose them into the following set of tensors with sample indices only:

$$\mathbb{E} \left[\widehat{\text{dd}_I H}_{i_0 i_1 i_2 i_3; \delta_0 \delta_1 \delta_2 \delta_3}^{(\ell)} \right] \equiv \frac{1}{n_{\ell-1}} \left[\delta_{i_0 i_1} \delta_{i_2 i_3} R_{\delta_0 \delta_1 \delta_2 \delta_3}^{(\ell)} + \delta_{i_0 i_2} \delta_{i_3 i_1} R_{\delta_0 \delta_2 \delta_3 \delta_1}^{(\ell)} + \delta_{i_0 i_3} \delta_{i_1 i_2} R_{\delta_0 \delta_3 \delta_1 \delta_2}^{(\ell)} \right], \quad (\infty.11)$$

$$\mathbb{E} \left[\widehat{\text{dd}_{II} H}_{i_1 i_2 i_3 i_4; \delta_1 \delta_2 \delta_3 \delta_4}^{(\ell)} \right] \equiv \frac{1}{n_{\ell-1}} \left[\delta_{i_1 i_2} \delta_{i_3 i_4} S_{\delta_1 \delta_2 \delta_3 \delta_4}^{(\ell)} + \delta_{i_1 i_3} \delta_{i_4 i_2} T_{\delta_1 \delta_3 \delta_4 \delta_2}^{(\ell)} + \delta_{i_1 i_4} \delta_{i_2 i_3} U_{\delta_1 \delta_4 \delta_2 \delta_3}^{(\ell)} \right]. \quad (\infty.12)$$

Thus, there are four new objects whose recursions we need to determine and whose scaling with depth we'll need to compute.

The ddNTKs both vanish in the first layer, just like the dNTK. Proceeding then from the second layer to deeper layers, after a bunch of tedious algebra and plenty of flipping around in the book to make various substitutions, you can find recursions for $R^{(\ell)}$, $S^{(\ell)}$, $T^{(\ell)}$, and $U^{(\ell)}$: (∞.177), (∞.179), (∞.180), and (∞.181), respectively. You can flip forward to take a look at them, but you probably won't want to... Importantly, these recursions in conjunction with the decompositions (∞.11) and (∞.12) altogether demonstrate that both ddNTKs have nontrivial order-1/ n statistics:

$$\mathbb{E} \left[\widehat{\text{dd}_I H}_{i_0 i_1 i_2 i_3; \delta_0 \delta_1 \delta_2 \delta_3}^{(\ell)} \right] = O\left(\frac{1}{n}\right), \quad \mathbb{E} \left[\widehat{\text{dd}_{II} H}_{i_1 i_2 i_3 i_4; \delta_1 \delta_2 \delta_3 \delta_4}^{(\ell)} \right] = O\left(\frac{1}{n}\right). \quad (\infty.13)$$

Thus, as we've been forecasting, we must include them in our finite-width analysis of training dynamics.

Focusing on the single-input statistics, we again make the layer-independent choices $C_b^{(\ell)} = C_b$, $C_W^{(\ell)} = C_W$, and $n_1 = \dots = n_{L-1} \equiv n$. Ignoring contributions that are subleading in $1/n$, in particular replacing the mean metric by the kernel $G^{(\ell)} \rightarrow K^{(\ell)}$ and the NTK mean by the frozen NTK $H^{(\ell)} \rightarrow \Theta^{(\ell)}$, the four recursions (∞.177), (∞.179), (∞.180), and (∞.181) together reduce to a form that at least fits on a page:

$$\begin{aligned} R^{(\ell+1)} &= \left(\chi_{\perp}^{(\ell)} \right)^2 R^{(\ell)} \\ &+ \lambda_W^{(\ell+1)} C_W \langle \sigma'' \sigma' \sigma' \sigma \rangle_{K^{(\ell)}} \left(\Theta^{(\ell)} \right)^2 + C_W^2 \langle \sigma''' \sigma' \sigma' \sigma' \rangle_{K^{(\ell)}} \left(\Theta^{(\ell)} \right)^3 \\ &+ \chi_{\perp}^{(\ell)} \left(\lambda_W^{(\ell+1)} \langle \sigma'' \sigma \rangle_{K^{(\ell)}} + C_W \Theta^{(\ell)} \langle \sigma''' \sigma' \rangle_{K^{(\ell)}} \right) \left(B^{(\ell)} + P^{(\ell)} \right) \\ &+ \chi_{\perp}^{(\ell)} \left(\lambda_W^{(\ell+1)} \langle \sigma' \sigma' \rangle_{K^{(\ell)}} + C_W \Theta^{(\ell)} \langle \sigma'' \sigma'' \rangle_{K^{(\ell)}} \right) P^{(\ell)}, \end{aligned} \quad (\infty.14)$$

$$S^{(\ell+1)} = \left(\chi_{\perp}^{(\ell)}\right)^2 S^{(\ell)} \quad (\infty.15)$$

$$+ C_W \lambda_W^{(\ell+1)} \langle \sigma' \sigma' \sigma' \sigma' \rangle_{K^{(\ell)}} \left(\Theta^{(\ell)}\right)^2 + C_W^2 \langle \sigma'' \sigma'' \sigma' \sigma' \rangle_{K^{(\ell)}} \left(\Theta^{(\ell)}\right)^3 \\ + \chi_{\perp}^{(\ell)} \left[\lambda_W^{(\ell+1)} \langle \sigma' \sigma' \rangle_{K^{(\ell)}} + C_W \Theta^{(\ell)} \langle \sigma'' \sigma'' \rangle_{K^{(\ell)}} \right] B^{(\ell)},$$

$$T^{(\ell+1)} = \left(\chi_{\perp}^{(\ell)}\right)^2 T^{(\ell)} \quad (\infty.16)$$

$$+ 2C_W \lambda_W^{(\ell+1)} \langle \sigma'' \sigma' \sigma' \sigma \rangle_{K^{(\ell)}} \left(\Theta^{(\ell)}\right)^2 + C_W^2 \langle \sigma'' \sigma'' \sigma' \sigma' \rangle_{K^{(\ell)}} \left(\Theta^{(\ell)}\right)^3 \\ + \left(\lambda_W^{(\ell+1)}\right)^2 \langle \sigma' \sigma' \sigma \sigma \rangle_{K^{(\ell)}} \Theta^{(\ell)} \\ + \left(\lambda_W^{(\ell+1)} \langle z \sigma' \sigma \rangle_{K^{(\ell)}} + C_W \Theta^{(\ell)} \langle z \sigma'' \sigma' \rangle_{K^{(\ell)}}\right)^2 \frac{F^{(\ell)}}{(K^{(\ell)})^2} \\ + 2\chi_{\perp}^{(\ell)} \left[\lambda_W^{(\ell+1)} (\langle \sigma'' \sigma \rangle_{K^{(\ell)}} + \langle \sigma' \sigma' \rangle_{K^{(\ell)}}) + C_W \Theta^{(\ell)} (\langle \sigma''' \sigma' \rangle_{K^{(\ell)}} + \langle \sigma'' \sigma'' \rangle_{K^{(\ell)}}) \right] Q^{(\ell)},$$

$$U^{(\ell+1)} = \left(\chi_{\perp}^{(\ell)}\right)^2 U^{(\ell)} + C_W^2 \langle \sigma'' \sigma'' \sigma' \sigma' \rangle_{K^{(\ell)}} \left(\Theta^{(\ell)}\right)^3. \quad (\infty.17)$$

Here, we also simplified these expressions by recalling the two susceptibilities, (5.50) and (5.51):

$$\chi_{\parallel}^{(\ell)} \equiv \frac{C_W}{K^{(\ell)}} \langle \sigma' \sigma z \rangle_{K^{(\ell)}}, \quad \chi_{\perp}^{(\ell)} \equiv C_W \langle \sigma' \sigma' \rangle_{K^{(\ell)}}. \quad (\infty.18)$$

ddNTK Scalings

Now, let's tune to criticality, and use our scaling ansatz (5.93) to find the critical exponents p_R , p_S , p_T , and p_U that describe the asymptotic depth scaling of $R^{(\ell)}$, $S^{(\ell)}$, $T^{(\ell)}$, and $U^{(\ell)}$, respectively.

To understand the relative size of these tensors controlling the means of the ddNTKs, we will again need to identify appropriate dimensionless ratios. Following the dimensional analysis logic from our dNTK discussion, cf. (11.63), let's look at our third-order update for preactivations (∞.4). Remembering again that we can only add terms that have the same dimensions, we see that

$$[z] = [\eta] [\epsilon] [\widehat{H}] = [\eta]^2 [\epsilon]^2 [\widehat{dH}] = [\eta]^3 [\epsilon]^3 [\widehat{dd_I H}] = [\eta]^3 [\epsilon]^3 [\widehat{dd_{II} H}]. \quad (\infty.19)$$

From the first equality, we see as before that $[\eta] [\epsilon] = [z] [\widehat{H}]^{-1}$, and so R , S , T , and U each have dimensions of NTK cubed:

$$[R] \equiv [\widehat{dd_I H}] = [\widehat{H}]^3 [z]^{-2}, \quad [S] = [T] = [U] \equiv [\widehat{dd_{II} H}] = [\widehat{H}]^3 [z]^{-2}. \quad (\infty.20)$$

This means that the proper dimensionless combinations for the ddNTKs are

$$\frac{R^{(\ell)} K^{(\ell)}}{n (\Theta^{(\ell)})^3} \sim \frac{1}{n} \left(\frac{1}{\ell}\right)^{p_R + p_0 - 3p_\Theta}, \quad \frac{S^{(\ell)} K^{(\ell)}}{n (\Theta^{(\ell)})^3} \sim \frac{1}{n} \left(\frac{1}{\ell}\right)^{p_S + p_0 - 3p_\Theta}, \quad (\infty.21)$$

$$\frac{T^{(\ell)} K^{(\ell)}}{n (\Theta^{(\ell)})^3} \sim \frac{1}{n} \left(\frac{1}{\ell}\right)^{p_T + p_0 - 3p_\Theta}, \quad \frac{U^{(\ell)} K^{(\ell)}}{n (\Theta^{(\ell)})^3} \sim \frac{1}{n} \left(\frac{1}{\ell}\right)^{p_U + p_0 - 3p_\Theta}, \quad (\infty.22)$$

where to get a normalized ratio, we multiplied by the kernel $K^{(\ell)}$ to account for the $[z]^{-2}$ and then divided by the three factors of the frozen NTK $\Theta^{(\ell)}$. Here, p_0 is the critical exponent for the kernel, and p_Θ is the critical exponent for the NTK.

Finally, as a brief aside, let us comment on one aspect of dimensionality that we've been ignoring but will soon become important. In particular, since the network output z is set to the true output y , they should really have the same dimensions:

$$[z] = [y]. \quad (\infty.23)$$

However, while the leading depth scaling of the preactivations is given by the kernel,

$$\mathbb{E} [z^{(\ell)} z^{(\ell)}] = K^{(\ell)} \sim \left(\frac{1}{\ell}\right)^{p_0}, \quad (\infty.24)$$

the true output y is fixed and doesn't scale with depth. When comparing the performance of networks of different depths, this suggests that it might be helpful to rescale the final network outputs as

$$z_{i;\delta} \rightarrow z_{i;\delta} \left(\frac{1}{L}\right)^{-p_0/2}, \quad (\infty.25)$$

effectively fixing the scaling of the overall network output as $p_0 = 0$, or equivalently rescale the training outputs as

$$y_{i;\tilde{\alpha}} \rightarrow y_{i;\tilde{\alpha}} \left(\frac{1}{L}\right)^{p_0/2}. \quad (\infty.26)$$

We will further see how this can affect the predictions of a fully-trained network on its test set in § $\infty.2.3$.¹

$K^* = 0$ Universality Class

For the $K^* = 0$ universality class, remember that we Taylor-expanded the activation function as

$$\sigma(z) = \sum_{p=0}^{\infty} \frac{\sigma_p}{p!} z^p, \quad (\infty.27)$$

¹For regression tasks, where we want to learn a vector of real numbers, this rescaling is appropriate. For classification tasks, where we want to learn a discrete probability distribution, we should rescale either the network outputs (before any softmax layer) or the raw output targets $y_{i;\delta}$ (again before any softmax layer) as in (10.37).

defined the following Taylor coefficient for convenience,

$$a_1 \equiv \left(\frac{\sigma_3}{\sigma_1} \right) + \frac{3}{4} \left(\frac{\sigma_2}{\sigma_1} \right)^2, \quad (\infty.28)$$

and required that all activation functions in this class satisfy $\sigma_0 = 0$ and $\sigma_1 \neq 0$.

To solve our single input ddNTK recursions, (∞.14)–(∞.17), you'll have to evaluate a few new Gaussian expectations, taking particular note of the fact that some of them now depend on the third derivative of the activation function. Finally, to tune to $K^* = 0$ criticality (5.90), we need to set the initialization hyperparameters as $C_b = 0$ and $C_W = 1/\sigma_1^2$; to implement the learning rate equivalence principle, we need to set our training hyperparameters as (9.95),

$$\lambda_b^{(\ell)} = \tilde{\lambda}_b \left(\frac{1}{\ell} \right)^{p_\perp} L^{p_\perp - 1}, \quad \lambda_W^{(\ell)} = \tilde{\lambda}_W \left(\frac{L}{\ell} \right)^{p_\perp - 1}. \quad (\infty.29)$$

For simplicity, let us also focus on odd activation functions, such as **tanh**, for which importantly $\sigma_2 = 0$ and $p_\perp = 1$.

Inspecting the single-input ddNTK recursions (∞.14)–(∞.17), we see that they depend on Gaussian expectations of preactivations as well as our previous solutions for all the other objects that we've considered: the NTK variance, the NTK–preactivation cross correlation, and the dNTK–preactivation cross correlation. Substituting in the solutions for all these quantities as needed – you'll have to flip around to find them, though most were reprinted in §11.3.2 for the dNTK analysis – we can find solutions to all the single-input ddNTK recursions:

$$R^{(\ell)} = -\frac{\ell^2}{48} \left[3\tilde{\lambda}_b + 4\frac{\tilde{\lambda}_W \sigma_1^2}{(-a_1)} \right] \left[\tilde{\lambda}_b + \frac{\tilde{\lambda}_W \sigma_1^2}{(-a_1)} \right]^2 (-a_1) + \cdots, \quad (\infty.30)$$

$$S^{(\ell)} = \frac{\ell^2}{12} \left[\tilde{\lambda}_b + \frac{\tilde{\lambda}_W \sigma_1^2}{(-a_1)} \right]^2 \tilde{\lambda}_W \sigma_1^2 + \cdots, \quad (\infty.31)$$

$$T^{(\ell)} = \frac{\ell^2}{32} \left[\tilde{\lambda}_b + \frac{\tilde{\lambda}_W \sigma_1^2}{(-a_1)} \right] (-a_1) \tilde{\lambda}_b^2 + \cdots, \quad (\infty.32)$$

$$U^{(\ell)} = \frac{1}{2} \left[\tilde{\lambda}_b + \frac{\tilde{\lambda}_W \sigma_1^2}{(-a_1)} \right]^3 (-a_1) + \cdots. \quad (\infty.33)$$

From these, we can read off the critical exponents as $p_R = p_S = p_T = -2$, and $p_U = 0$, and we see that the dimensionless ratios are given by

$$\frac{R^{(\ell)} K^{(\ell)}}{n(\Theta^{(\ell)})^3} = -\frac{1}{48} \left[\frac{3\tilde{\lambda}_b + 4\tilde{\lambda}_W \sigma_1^2 / (-a_1)}{\tilde{\lambda}_b + \tilde{\lambda}_W \sigma_1^2 / (-a_1)} \right] \frac{\ell}{n} + \cdots, \quad (\infty.34)$$

$$\frac{S^{(\ell)} K^{(\ell)}}{n(\Theta^{(\ell)})^3} = \frac{1}{12} \left[\frac{\tilde{\lambda}_W \sigma_1^2 / (-a_1)}{\tilde{\lambda}_b + \tilde{\lambda}_W \sigma_1^2 / (-a_1)} \right] \frac{\ell}{n} + \cdots, \quad (\infty.35)$$

$$\frac{T^{(\ell)}K^{(\ell)}}{n\left(\Theta^{(\ell)}\right)^3}=\frac{1}{32}\left[\frac{\tilde{\lambda}_b}{\tilde{\lambda}_b+\tilde{\lambda}_W\sigma_1^2/(-a_1)}\right]^2\frac{\ell}{n}+\cdots, \tag{\infty.36}$$

$$\frac{U^{(\ell)}K^{(\ell)}}{n\left(\Theta^{(\ell)}\right)^3}=\frac{1}{2}\frac{1}{\ell n}+\cdots. \tag{\infty.37}$$

This means that $R^{(\ell)}$, $S^{(\ell)}$, and $T^{(\ell)}$ all scale according to our leading effective theory cutoff as

$$p_R+p_0-3p_\Theta=-1 \qquad p_S+p_0-3p_\Theta=-1, \qquad p_T+p_0-3p_\Theta=-1. \tag{\infty.38}$$

Thus, we see that the leading-order finite-width dynamics of $K^\star=0$ activation functions have contributions from the first ddNTK, $\widehat{\text{dd}_\text{I}H}$, via $R^{(\ell)}$, and from the second ddNTK, $\widehat{\text{dd}_\text{II}H}$, via $S^{(\ell)}$ and $T^{(\ell)}$.²

Scale-Invariant Universality Class

Perhaps the most important fact to remember about nonlinear scale-invariant activation functions (2.13),

$$\sigma(z)=\begin{cases} a_+z, & z\geq 0, \\ a_-z, & z<0, \end{cases} \tag{\infty.39}$$

with $a_+\neq -a_-$, is that they are not smooth: their first derivative is a step function centered at the origin; their second derivative is a Dirac delta function, (2.32),

$$\sigma'(z)=\begin{cases} a_+, & z\geq 0, \\ a_-, & z<0, \end{cases}, \qquad \sigma''(z)=(a_+-a_-)\delta(z); \tag{\infty.40}$$

and higher derivatives will involve derivatives of the Dirac delta function. Inspecting again the single-input ddNTK recursions (∞.14)–(∞.17), the kink in these activation functions and the presence of Gaussian expectations with up to three derivatives of $\sigma(z)$ should scare you, especially if you heeded our warning at the end of §5.5. In fact, if you formally try to evaluate some of these expectations, particularly $\langle\sigma''\sigma''\rangle_K$ and others related to it via integration by parts, you’ll find that they want to blow up, even if you use all the *magic tricks* from physics that you might have at your disposal for trying to make sense of divergent integrals.³

The divergence of these Gaussian correlators is actually telling us that there is something very wrong with our expansions for the updates to the preactivations (∞.4),

²If we relaxed the restriction for the activation function to be odd, we’d find the same scalings for $R^{(\ell)}$, $S^{(\ell)}$, and $T^{(\ell)}$ – though with different coefficients – and we’d find that $U^{(\ell)}$ was up by a factor of ℓ , but still subleading overall.

³We were somewhat lucky in §11.3.1 when analyzing the dNTK: all the higher-derivative Gaussian integrals that we needed simplified via integration by parts and gave finite – in fact, vanishing – answers, cf. (11.67) and (11.69).

the NTK (∞.6), and the dNTK (∞.9). In particular, the Taylor expansion in the global learning rate η breaks down for these nonsmooth activation functions and doesn't accurately describe how a network is updated. As a result, our approach for solving the finite-width dynamics will not work for nonlinear scale-invariant activation functions.

To understand why, let's consider an extremely simple model of a network, a single neuron with a bias:

$$z(x) = \sigma(x + b). \quad (\infty.41)$$

Here, the input x and the output z are both scalars, and for the activation function $\sigma(z)$ we'll pick the ReLU, with $a_+ = 1$ and $a_- = 0$. Accordingly, for a particular input x such that $x + b > 0$, the activation fires, and the output is $z(x) = x + b$; for a particular input x' such that $x' + b < 0$, the activation doesn't fire, and the output is $z(x') = 0$.

Now, let's consider a gradient-descent update to the parameters with a training example $x > -b$ such that the activation fires. Then, the bias updates as

$$\vec{db} = -\eta \frac{dz}{db} \epsilon = -\eta \epsilon = O(\eta), \quad (\infty.42)$$

where ϵ is the error factor of the loss, depending on the true output y . Now, if $x + b + \vec{db} > 0$, then the change in the output is

$$\vec{dz} = -\eta \epsilon = O(\eta) \quad (\infty.43)$$

and would be perfectly described by our Taylor expansion (∞.4). However, if the training error is large enough such that $x + b + \vec{db} = x + b - \eta \epsilon < 0$, then the activation turns off and

$$\vec{dz} = -x - b = O(1). \quad (\infty.44)$$

Importantly, this update is independent of our expansion parameter η , and a Taylor expansion in η cannot detect this discontinuity at $\eta = (x + b)/\epsilon$.

Thus, for the ReLU, any time an activation crosses its firing threshold, it can contribute to the gradient-descent update in an η -independent way. Empirically, if you try to measure the NTK update for a deep MLP consisting of ReLU activation functions, you don't find anything consistent with the expansion (∞.6). Instead, for the appropriately normalized quantity, you'll find an $\sim 1/\sqrt{n}$ scaling with width and a linear scaling with depth ℓ , in contrast to the ℓ/n scaling expected from our perturbative formalism.⁴ From this we reach the unfortunate conclusion that we'll have to give up on describing the finite-width dynamics of nonlinear scale-invariant activation functions using these methods.

Note that everything we have discussed for nonlinear scale-invariant activation functions with respect to criticality and the infinite-width training dynamics is perfectly

⁴It's tempting to think that this $\sim 1/\sqrt{n}$ scaling arises from accumulating the probabilities that one of the Ln total hidden-layer $n \gg 1$ activations experiences an η -independent $O(1)$ change after the gradient-descent step.

fine and experimentally validated, and the presence of a nonzero dNTK at finite width is still indicative of a dynamical NTK and representation learning at finite width. The problem we just described only affects the analysis we're going to perform in the following section for the training dynamics of finite-width networks, and the breakdown of the Taylor expansion just means that we will be unable to give a quantitative picture of representation learning for these nonsmooth activation functions.⁵ So, if you do want to understand this, you'll probably need an entirely new approach.

This leaves us one activation function in the entire scale-invariant universality class: the **linear** activation function used for deep linear networks. Tuning to criticality, $C_b = 0$ and $C_W = 1$, which fixes $\chi = 1$, and choosing layer-independent learning rates (9.94) as

$$\lambda_b^{(\ell)} = \frac{\tilde{\lambda}_b}{L}, \quad \lambda_W^{(\ell)} = \frac{\tilde{\lambda}_W}{L}, \quad (\infty.45)$$

we can solve the single-input ddNTK recursions ($\infty.14$)–($\infty.17$). Note that for every term in the U -recursion, ($\infty.17$), and for every term but one in the R -recursion, ($\infty.14$), the Gaussian expectations of activations involve second derivatives or third derivatives, which vanish for a **linear** function. For the one term in the R -recursion that does not vanish, it is also multiplied by $P^{(\ell)}$, which vanishes for all scale-invariant activations, cf. (11.74). This means that

$$R^{(\ell)} = 0, \quad U^{(\ell)} = 0 \quad (\infty.46)$$

and in particular that the first ddNTK, $\widehat{\text{dd}_1 H}$, does not contribute to the deep linear network's dynamics at $O(1/n)$ since it's entirely determined by $R^{(\ell)}$. However, for the S -recursion, ($\infty.15$), and the T -recursion, ($\infty.16$), we find something nontrivial: these recursions can be exactly solved as

$$S^{(\ell)} = \frac{\ell^2(\ell^2 - 1)}{12L^3} (\tilde{\lambda}_b + \tilde{\lambda}_W K^*)^2 \tilde{\lambda}_W, \quad (\infty.47)$$

$$T^{(\ell)} = \frac{\ell^2(\ell^2 - 1)}{12L^3} (\tilde{\lambda}_b + \tilde{\lambda}_W K^*) \tilde{\lambda}_W^2 K^*, \quad (\infty.48)$$

from which we see that the critical exponents are $p_S = p_T = -4$. Finally, the dimensionless ratios are

$$\frac{S^{(\ell)} K^{(\ell)}}{n (\Theta^{(\ell)})^3} = \frac{1}{12} \left[\frac{\tilde{\lambda}_W K^*}{\tilde{\lambda}_b + \tilde{\lambda}_W K^*} \right] \frac{\ell}{n} + \dots, \quad (\infty.49)$$

$$\frac{T^{(\ell)} K^{(\ell)}}{n (\Theta^{(\ell)})^3} = \frac{1}{12} \left[\frac{\tilde{\lambda}_W K^*}{\tilde{\lambda}_b + \tilde{\lambda}_W K^*} \right]^2 \frac{\ell}{n} + \dots, \quad (\infty.50)$$

⁵However, our analysis applies to any of the smoothed versions of the **ReLU** that permit a type of criticality, cf. our criticality discussion of the **SWISH** and **GELU** in §5.3.4. In particular, the training dynamics we'll work out in the next section describe these networks. These dynamical solutions, in conjunction with the output-layer solutions to all their associated recursions, will accurately characterize such fully-trained **ReLU**-like networks in practice.

and we see that these scale according to our leading effective theory cutoff as

$$p_S + p_0 - 3p_\Theta = -1, \quad p_T + p_0 - 3p_\Theta = -1. \quad (\infty.51)$$

In conclusion, we see that the second ddNTK, $\widehat{\text{dd}}_{\text{II}}H$, contributes via $S^{(\ell)}$ and $T^{(\ell)}$ to the dynamics of deep linear networks at leading order.

∞.2 Training at Finite Width

Now that we understand the joint statistics of the preactivations, the NTK, the dNTK, and the ddNTKs, we have nearly all the tools we need in order to evaluate the distribution of *fully-trained* networks at finite width and nonzero depth. To see why, recall our finite-width expansion of the network output evolution (∞.4):

$$\begin{aligned} z_{i;\delta}^{(L)}(t=1) &= z_{i;\delta}^{(L)} - \eta \sum_{j,\tilde{\alpha}} \widehat{H}_{ij;\delta\tilde{\alpha}}^{(L)} \epsilon_{j;\tilde{\alpha}}^{(L)} + \frac{\eta^2}{2} \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2} \widehat{\text{d}H}_{ij_1j_2;\delta\tilde{\alpha}_1\tilde{\alpha}_2}^{(L)} \epsilon_{j_1;\tilde{\alpha}_1}^{(L)} \epsilon_{j_2;\tilde{\alpha}_2}^{(L)} \\ &\quad - \frac{\eta^3}{6} \sum_{j_1,j_2,j_3,\tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3} \widehat{\text{dd}}_{\text{I}}H_{ij_1j_2j_3;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3}^{(L)} \epsilon_{j_1;\tilde{\alpha}_1}^{(\ell)} \epsilon_{j_2;\tilde{\alpha}_2}^{(\ell)} \epsilon_{j_3;\tilde{\alpha}_3}^{(\ell)} + O\left(\frac{1}{n^2}\right). \end{aligned} \quad (\infty.52)$$

Importantly, all the quantities on the right-hand side of the network update (∞.52) – $z_{i;\delta}^{(L)}$, $\widehat{H}_{ij;\delta\tilde{\alpha}}^{(L)}$, $\epsilon_{j;\tilde{\alpha}}^{(L)}$, $\widehat{\text{d}H}_{ij_1j_2;\delta\tilde{\alpha}_1\tilde{\alpha}_2}^{(L)}$, and $\widehat{\text{dd}}_{\text{I}}H_{ij_1j_2j_3;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3}^{(L)}$ – are evaluated at initialization and thus are determined completely by the statistics of the joint preactivation–NTK–dNTK–ddNTKs distribution,

$$p\left(z^{(L)}, \widehat{H}^{(L)}, \widehat{\text{d}H}^{(L)}, \widehat{\text{dd}}_{\text{I}}H^{(L)}, \widehat{\text{dd}}_{\text{II}}H^{(L)} \middle| \mathcal{D}\right), \quad (\infty.53)$$

that we spent the majority of this book evaluating in detail. Accordingly – just as we did before at infinite width (§10) – we could use the joint distribution (∞.53) to compute the statistics of the fully-trained network outputs after the update (∞.52), *if* we tuned a single step of gradient descent for each realization of a network so that we landed on the minimum of the training loss $z_{i;\tilde{\alpha}}^{(L)}(t=1) = y_{i;\tilde{\alpha}}$.

More generally, *if* we trained the network with T steps of gradient descent such that

$$z_{i;\tilde{\alpha}}^{(L)}(t=T) = y_{i;\tilde{\alpha}}, \quad \text{for all } \tilde{\alpha} \in \mathcal{A}, \quad (\infty.54)$$

and *if* we determined how to express the output of such a fully-trained network for general inputs $\delta \in \mathcal{D}$ as a functional of our statistical variables at initialization,

$$z_{i;\delta}^{(L)}(T) \equiv \left[z_{i;\delta}^{(L)}(t=T) \right] \left(z^{(L)}, \widehat{H}^{(L)}, \widehat{\text{d}H}^{(L)}, \widehat{\text{dd}}_{\text{I}}H^{(L)}, \widehat{\text{dd}}_{\text{II}}H^{(L)} \right), \quad (\infty.55)$$

then we could give an analytical expression for the distribution of *fully-trained network outputs*:

$$p\left(z^{(L)}(T)\right). \quad (\infty.56)$$

The equation $(\infty.54)$ is our *fully-trained condition*, and the distribution $(\infty.56)$ completely describes the ensemble of finite-width networks at the end of training. The theoretical understanding of this distribution is exactly the goal we set for ourselves at the beginning of this book in §0. The only thing left for us to work out is the functional $(\infty.55)$; to do so, we first need to figure out what kind of steps to take in order to fully train these finite-width networks.

Now, recall from §10.2.2 that in the infinite-width limit, the fully-trained network solution $(\infty.55)$ had an *algorithm independence*: the distribution at the end of training didn't depend on the details of the optimization algorithm, and thus we could perform our theoretical analysis with any algorithm we wanted. In contrast, at finite width, the fully-trained solution $(\infty.55)$ will have an **algorithm dependence**: different fully-trained solutions will make different test-set predictions depending on the details of the particular optimization algorithm used to train the network, even when holding fixed the statistical variables at initialization and their associated initialization and training hyperparameters. Encouragingly, the solutions will nonetheless take a universal form, with the nonuniversal details of the particular training algorithm captured by six projective tensors: cf. $(\infty.154)$. Thus, we will be able to very generally study the distribution of fully-trained networks at finite width by working out such solutions.

With that in mind, in this section we'll present fully-trained solutions for two different optimization algorithms. First, in § $\infty.2.1$ we'll take *two* Newton-like steps in order to satisfy our fully-trained condition $(\infty.54)$. While practically infeasible for large training sets, this training algorithm is rich in pedagogical value, emphasizing the way in which a finite-width network needs to adapt its representation to the training data in order to minimize its training error. Then, in § $\infty.2.2$ we'll analytically solve the dynamics of the vanilla gradient descent at order $1/n$ and obtain a slightly different ensemble of fully-trained finite-width networks. This algorithm is not only practically implementable but also quite often used to optimize real neural networks, and our corresponding solution is an actual theoretical description of such fully-trained networks. Together, these solutions will help us understand the ways in which the details of the optimization algorithm can affect the corresponding fully-trained solution. Finally, in § $\infty.2.3$ we'll be able to generally analyze the predictions of these different fully-trained networks on novel examples from the test set.

Throughout this section, we will declutter the notation a bit by dropping the layer indices, since to understand training we only need to focus on the network output at layer $\ell = L$.

An Infinite-Width Giant Leap at Finite Width

Before we begin, let's first review the giant leap that we took in §10.2 at infinite width. From the finer-grained perspective of finite width, we'll see that our leap actually missed the minimum, exhibiting training errors of order $1/n$. However, our new eyes on this leap will be instructive, as they will help us see how we can correct for these errors and reduce the finite-width training error even further.

Recall from §10.2 that, in order to fully train an infinite-width network in a single step, we needed to make a *second-order update* of the form

$$\bar{d}\theta_\mu = - \sum_{\nu, \tilde{\alpha}_1, \tilde{\alpha}_2, i} \eta \lambda_{\mu\nu} \kappa^{\tilde{\alpha}_1 \tilde{\alpha}_2} \frac{dz_{i; \tilde{\alpha}_1}}{d\theta_\nu} (z_{i; \tilde{\alpha}_2} - y_{i; \tilde{\alpha}_2}), \quad (\infty.57)$$

which we interpreted either (i) as a *generalized training algorithm* (10.19) optimizing the standard MSE loss (10.5), or (ii) as a standard (tensorial) gradient-descent step (7.11) optimizing a *generalized MSE loss* (10.22). Here, $\kappa^{\tilde{\alpha}_1 \tilde{\alpha}_2}$ was called the *Newton tensor*, and – in the first understanding of (∞.57) – could be interpreted as allowing us to take anisotropic steps in training sample space.

With this type of parameter update, a finite-width network output will evolve as

$$\begin{aligned} z_{i; \delta}(t=1) &= z_{i; \delta} - \eta \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\delta \tilde{\alpha}_1} \kappa^{\tilde{\alpha}_1 \tilde{\alpha}_2} (z_{i; \tilde{\alpha}_2} - y_{i; \tilde{\alpha}_2}) - \eta \sum_{j, \tilde{\alpha}_1, \tilde{\alpha}_2} \widehat{\Delta H}_{ij; \delta \tilde{\alpha}_1} \kappa^{\tilde{\alpha}_1 \tilde{\alpha}_2} (z_{j; \tilde{\alpha}_2} - y_{j; \tilde{\alpha}_2}) \\ &\quad + \frac{\eta^2}{2} \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \widehat{dH}_{ij_1 j_2; \delta \tilde{\alpha}_1 \tilde{\alpha}_2} \kappa^{\tilde{\alpha}_1 \tilde{\alpha}_3} \kappa^{\tilde{\alpha}_2 \tilde{\alpha}_4} (z_{j_1; \tilde{\alpha}_3} - y_{j_1; \tilde{\alpha}_3}) (z_{j_2; \tilde{\alpha}_4} - y_{j_2; \tilde{\alpha}_4}) \\ &\quad - \frac{\eta^3}{6} \sum_{j_1, j_2, j_3, \tilde{\alpha}_1, \dots, \tilde{\alpha}_6} \widehat{dd_I H}_{ij_1 j_2 j_3; \delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \kappa^{\tilde{\alpha}_1 \tilde{\alpha}_4} \kappa^{\tilde{\alpha}_2 \tilde{\alpha}_5} \kappa^{\tilde{\alpha}_3 \tilde{\alpha}_6} \\ &\quad \times (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) + O\left(\frac{1}{n^2}\right), \end{aligned} \quad (\infty.58)$$

where here we've made our usual decomposition of the NTK into a mean and fluctuation:

$$\widehat{H}_{ij; \delta \tilde{\alpha}} \equiv \delta_{ij} H_{\delta \tilde{\alpha}} + \widehat{\Delta H}_{ij; \delta \tilde{\alpha}}. \quad (\infty.59)$$

In terms of such an update, our fully-trained condition (∞.54) after a single step $T = 1$,

$$z_{i; \tilde{\alpha}}(t=1) = y_{i; \tilde{\alpha}}, \quad \text{for all } \tilde{\alpha} \in \mathcal{A}, \quad (\infty.60)$$

can be written as

$$\begin{aligned} z_{i; \tilde{\alpha}} - y_{i; \tilde{\alpha}} &= \eta \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\tilde{\alpha} \tilde{\alpha}_1} \kappa^{\tilde{\alpha}_1 \tilde{\alpha}_2} (z_{i; \tilde{\alpha}_2} - y_{i; \tilde{\alpha}_2}) + \eta \sum_{j, \tilde{\alpha}_1, \tilde{\alpha}_2} \widehat{\Delta H}_{ij; \tilde{\alpha} \tilde{\alpha}_1} \kappa^{\tilde{\alpha}_1 \tilde{\alpha}_2} (z_{j; \tilde{\alpha}_2} - y_{j; \tilde{\alpha}_2}) \\ &\quad - \frac{\eta^2}{2} \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \widehat{dH}_{ij_1 j_2; \tilde{\alpha} \tilde{\alpha}_1 \tilde{\alpha}_2} \kappa^{\tilde{\alpha}_1 \tilde{\alpha}_3} \kappa^{\tilde{\alpha}_2 \tilde{\alpha}_4} (z_{j_1; \tilde{\alpha}_3} - y_{j_1; \tilde{\alpha}_3}) (z_{j_2; \tilde{\alpha}_4} - y_{j_2; \tilde{\alpha}_4}) \\ &\quad + \frac{\eta^3}{6} \sum_{j_1, j_2, j_3, \tilde{\alpha}_1, \dots, \tilde{\alpha}_6} \widehat{dd_I H}_{ij_1 j_2 j_3; \tilde{\alpha} \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \kappa^{\tilde{\alpha}_1 \tilde{\alpha}_4} \kappa^{\tilde{\alpha}_2 \tilde{\alpha}_5} \kappa^{\tilde{\alpha}_3 \tilde{\alpha}_6} \\ &\quad \times (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) + O\left(\frac{1}{n^2}\right). \end{aligned} \quad (\infty.61)$$

This new giant-leap condition (∞.61) perhaps seems a little daunting, and further it's not obvious that there's any particular choice of Newton tensor $\kappa^{\tilde{\alpha}_1\tilde{\alpha}_2}$ that can land us on the minimum. Nonetheless, we do expect that our infinite-width solution should be near the true finite-width solution, up to errors of order $O(1/n)$.⁶

With that in mind, as a first step let's try our infinite-width giant leap (10.27) and see where we land. This infinite-width giant leap had an interpretation as *Newton's method* and was given by the second-order update (∞.57), with a particular choice of the product of the global learning rate and the Newton tensor,

$$\eta\kappa^{\tilde{\alpha}_1\tilde{\alpha}_2} = \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2}, \quad (\infty.62)$$

where the *inverse* NTK mean submatrix $\tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2}$ was defined implicitly via

$$\sum_{\tilde{\alpha}_2 \in \mathcal{A}} \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} \tilde{H}_{\tilde{\alpha}_2\tilde{\alpha}_3} = \delta^{\tilde{\alpha}_1}_{\tilde{\alpha}_3}. \quad (\infty.63)$$

As always, the tilde on $\tilde{H}_{\tilde{\alpha}_1\tilde{\alpha}_2}$ emphasizes that it's an $N_{\mathcal{A}} \times N_{\mathcal{A}}$ -dimensional submatrix of the NTK mean evaluated on pairs of training inputs *only*. By now this distinction should be familiar enough that we will stop belaboring it.

More importantly, here we've used the inverse of the full NTK mean $\tilde{H}_{\tilde{\alpha}_1\tilde{\alpha}_2}$ rather than the infinite-width frozen NTK $\tilde{\Theta}_{\tilde{\alpha}_1\tilde{\alpha}_2}$. To explain why, let us recall from (9.3) that the NTK mean receives a series of corrections at each order in the $1/n$ expansion, of which the leading-order piece is the frozen NTK (9.4). Since we're now working at order $1/n$, we should in particular take into account the next-to-leading-order (NLO) $1/n$ correction to the NTK mean $H_{\tilde{\alpha}_1\tilde{\alpha}_2}^{\{1\}}$ by working with $\tilde{H}_{\tilde{\alpha}_1\tilde{\alpha}_2}$ instead of $\tilde{\Theta}_{\tilde{\alpha}_1\tilde{\alpha}_2}$.⁷

Substituting our infinite-width giant-leap Newton tensor (∞.62) into our giant-leap condition (∞.61) at finite width and rearranging, we get

$$\begin{aligned} 0 = & \sum_{j, \tilde{\alpha}_1, \tilde{\alpha}_2} \widehat{\Delta H}_{ij; \tilde{\alpha}\tilde{\alpha}_1} \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (z_{j; \tilde{\alpha}_1} - y_{j; \tilde{\alpha}_1}) \\ & - \frac{1}{2} \sum_{\substack{j_1, j_2, \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_4}} \widehat{dH}_{ij_1j_2; \tilde{\alpha}\tilde{\alpha}_1\tilde{\alpha}_2} \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_2\tilde{\alpha}_4} (z_{j_1; \tilde{\alpha}_3} - y_{j_1; \tilde{\alpha}_3})(z_{j_2; \tilde{\alpha}_4} - y_{j_2; \tilde{\alpha}_4}) \end{aligned} \quad (\infty.64)$$

⁶For deep networks, we know that technically the corrections will be of order $O(L/n)$, corresponding to the cutoff scale of our effective theory description.

⁷While we won't show it explicitly, we expect that this *NLO NTK mean*, $H_{\tilde{\alpha}_1\tilde{\alpha}_2}^{\{1\}(\ell)}$, will have a solution that scales like $O(1/n)$ as compared to the frozen NTK; this would be analogous to what we found for the NLO metric $G_{\tilde{\alpha}_1\tilde{\alpha}_2}^{\{1\}(\ell)}$, cf. the discussion in §5.4 after (5.143). In particular, if we make a $1/n$ expansion for our training hyperparameters $\lambda_b^{(\ell)}$ and $\lambda_W^{(\ell)}$ as we did for our initialization hyperparameters in (5.138) and (5.139), then we will have extra freedom in the subleading hyperparameters $\lambda_b^{(\ell)\{1\}}$ and $\lambda_W^{(\ell)\{1\}}$ to eliminate the growing-in- ℓ contribution to $H_{\tilde{\alpha}_1\tilde{\alpha}_2}^{\{1\}(\ell)}$. Overall, this will make the NLO correction to the NTK mean scale as $O(1/n)$, subleading to the leading finite-width effects which scale as $O(L/n)$: in the language of RG flow, the NLO NTK mean is *marginal*. In practice, such a contribution is negligible and can thus be neglected for networks of any real depth.

$$\begin{aligned}
 & + \frac{1}{6} \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_6}} \widehat{\text{dd}_1 H}_{i_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_4} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_6} \\
 & \times (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) + O\left(\frac{1}{n^2}\right).
 \end{aligned}$$

Thus, we actually missed the minimum: for the network to be fully trained, the right-hand side of (∞.64) should have vanished, while here it is clearly nonzero in general. Taking a step back, it's clear now that what we actually found at infinite width in §10.2 was

$$z_{i; \tilde{\alpha}}(t = 1) - y_{i; \tilde{\alpha}} = O\left(\frac{1}{n}\right). \quad (\infty.65)$$

In other words, our networks were fully trained only at leading order, and (∞.64) gives an explicit expression for the $1/n$ correction.

Disentangling a little further, there are two such corrections at order $1/n$: the first correction – the first term on the right-hand side of (∞.64) – arises from the instantiation-to-instantiation fluctuations of the NTK across different realizations of the biases and weights at initialization; the second correction – comprised of the second and third terms on the right-hand side of (∞.64) – arises from nonlinear changes in the output as we take our step. In particular, this second correction is a *bona fide* manifestation of representation learning at finite width, accounting for the fact that the network's *effective features* are evolving. If we properly account for these two types of $1/n$ corrections, we should be able to attain a training error of order $1/n^2$:

$$z_{i; \tilde{\alpha}}(T) - y_{i; \tilde{\alpha}} = O\left(\frac{1}{n^2}\right). \quad (\infty.66)$$

That is, we should be able to improve our effective theory of fully-trained networks, quantitatively by another multiplicative factor of $1/n$ and qualitatively by properly including representation learning into such an effective description.

Our first approach (§∞.2.1) to attain such effectively-zero training error, (∞.66), is to continue to engineer theoretical giant leaps so as to account for both the instantiation-to-instantiation fluctuations of the NTK and the effect of the dNTK and the first ddNTK.

Another approach (§∞.2.2) is to simply use the vanilla tensorial gradient descent algorithm as we do in practice; in that case, we will have to not only account for the dynamics of the network output but also account for the dynamics of the NTK and dNTK. After doing so, we will see that we can iteratively decrease the training loss to zero after many many such steps.

∞.2.1 A Small Step Following a Giant Leap

Here we'll train our networks with *two* second-order updates. For the first update, a giant leap, we'll need to further generalize our theoretical optimization algorithm in order to properly account for the instantiation-to-instantiation fluctuations of the NTK. In the

second update, a small step, we'll be able to account for the $1/n$ change in representation due to the nonzero dNTK and ddNTK, ultimately landing on the minimum of the loss as $(\infty.66)$. In particular, we can think of these updates as loosely corresponding to distinct phases of training that arise when implementing gradient-based training of neural networks in practice.

First Update: One Final Generalization of Gradient Decent

Please flip back to take a closer look at our unsatisfied condition for fully training our networks $(\infty.64)$. Right away, you should notice a serious problem in satisfying this constraint: the NTK fluctuation, $\widehat{\Delta H}_{ij;\tilde{\alpha}\tilde{\alpha}_1}$, the dNTK, $\widehat{dH}_{ij_1j_2;\tilde{\alpha}\tilde{\alpha}_1\tilde{\alpha}_2}$, and the ddNTK, $\widehat{dd_1H}_{ij_1j_2j_3;\tilde{\alpha}\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3}$, all mix different output components together in an update; i.e., the j -th component of the prediction error $z_{j;\tilde{\alpha}} - y_{j;\tilde{\alpha}}$ at initialization affects the i -th component of the output $z_{i;\tilde{\alpha}}$ after the update, even for $i \neq j$. This stands in contrast to what we found for an infinite-width update in §10.1.1, where there was no such mixing or *wiring* of output components. Meanwhile, we did see a similar wiring effect for Bayesian inference at finite width, as we discussed in depth in §6.4.2.

While such wiring at finite width is fantastic from a practitioner's standpoint, it makes our theoretical work slightly more complicated. In particular, in order to satisfy the training constraint $(\infty.64)$, we will need to further generalize our second-order update $(\infty.57)$. Following in the footsteps of our previous two generalizations, (7.11) and (10.18) , let's make one final generalization,

$$\eta \rightarrow \eta \lambda_{\mu\nu} \rightarrow \eta \lambda_{\mu\nu} \kappa_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2} \rightarrow \eta \lambda_{\mu\nu} \kappa_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2}, \quad (\infty.67)$$

with the final form of our theoretical update given by

$$d\theta_\mu = - \sum_{\nu, \tilde{\alpha}_1, \tilde{\alpha}_2, i, j} \eta \lambda_{\mu\nu} \kappa_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2} \frac{dz_{i;\tilde{\alpha}_1}}{d\theta_\nu} \epsilon_{j;\tilde{\alpha}_2}. \quad (\infty.68)$$

Here, we also introduced a further *generalized* Newton tensor $\kappa_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2}$ with output-component indices. This flexibility will allow us to resolve the mixing of the output components in the residual training error from our infinite-width leap $(\infty.64)$.⁸

⁸Similarly to the generalized MSE loss (10.22) discussed in §10.2.1, we can alternatively think of this further-generalized second-order update $(\infty.68)$ optimizing the standard MSE loss as instead arising from a standard first-order (tensorial) gradient descent update (7.11) on a further-generalized MSE loss,

$$\mathcal{L}_{\mathcal{A}, \kappa}(\theta) \equiv \frac{1}{2} \sum_{i_1, i_2=1}^{n_L} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} \kappa_{i_1 i_2}^{\tilde{\alpha}_1 \tilde{\alpha}_2} (z_{i_1; \tilde{\alpha}_1} - y_{i_1; \tilde{\alpha}_1}) (z_{i_2; \tilde{\alpha}_2} - y_{i_2; \tilde{\alpha}_2}), \quad (\infty.69)$$

as long as the Newton tensor $\kappa_{i_1 i_2}^{\tilde{\alpha}_1 \tilde{\alpha}_2}$ is *also* assumed to be symmetric under the exchange of paired indices $(i_1, \tilde{\alpha}_1) \leftrightarrow (i_2, \tilde{\alpha}_2)$. This symmetry will in fact be present for both our first update, the giant leap, and our second update, the small step. With this interpretation $(\infty.69)$, our generalized Newton tensor $\kappa_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2}$ acts as a metric on sample space, through its sample indices $\tilde{\alpha}_1, \tilde{\alpha}_2$, and on output-component space, through its L -th-layer neural indices i, j .

Note importantly that the μ, ν indices of $\lambda_{\mu\nu}$ are very different from the i, j indices of $\kappa_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2}$: the former each runs over all P parameters, while the latter each only runs over the n_L components of the network output. In particular, while learning-rate tensor $\lambda_{\mu\nu}$ lets us control how the gradient of the ν -th parameter affects the update to the μ -th parameter, the i, j indices of the generalized Newton tensor $\kappa_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2}$ instead control how the network's *features* (10.137) $dz_{i;\tilde{\alpha}_1}/d\theta_\nu$ are combined with the error factor $\epsilon_{j;\tilde{\alpha}_2}$ in order to make the update.⁹ Allowing for a $\kappa_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2}$ with nonzero off-diagonal components in the i, j indices, we can precisely tune the *wiring* or mixing of output components that occurs in a particular update.

Finally, plugging our final-form second-order update (∞.68) back into the $1/n$ expansion of the network update (∞.4) and using the NTK, dNTK, and first ddNTK definitions, we can see how the network output changes after making an update with this new optimization algorithm:

$$\begin{aligned} z_{i;\delta}(t=1) &= z_{i;\delta} - \eta \sum_{j,\tilde{\alpha}_1,\tilde{\alpha}_2} H_{\delta\tilde{\alpha}_1} \kappa_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2} (z_{j;\tilde{\alpha}_2} - y_{j;\tilde{\alpha}_2}) - \eta \sum_{j,k,\tilde{\alpha}_1,\tilde{\alpha}_2} \widehat{\Delta H}_{ij;\delta\tilde{\alpha}_1} \kappa_{jk}^{\tilde{\alpha}_1 \tilde{\alpha}_2} (z_{k;\tilde{\alpha}_2} - y_{k;\tilde{\alpha}_2}) \\ &+ \frac{\eta^2}{2} \sum_{\substack{j_1,j_2,k_1,k_2, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4}} \widehat{\text{d}H}_{ij_1j_2;\delta\tilde{\alpha}_1\tilde{\alpha}_2} \kappa_{j_1k_1}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \kappa_{j_2k_2}^{\tilde{\alpha}_2 \tilde{\alpha}_4} (z_{k_1;\tilde{\alpha}_3} - y_{k_1;\tilde{\alpha}_3}) (z_{k_2;\tilde{\alpha}_4} - y_{k_2;\tilde{\alpha}_4}) \\ &- \frac{\eta^3}{6} \sum_{\substack{j_1,j_2,j_3,k_1,k_2,k_3 \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \widehat{\text{dd}_I H}_{ij_1j_2j_3;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} \kappa_{j_1k_1}^{\tilde{\alpha}_1 \tilde{\alpha}_4} \kappa_{j_2k_2}^{\tilde{\alpha}_2 \tilde{\alpha}_5} \kappa_{j_3k_3}^{\tilde{\alpha}_3 \tilde{\alpha}_6} \\ &\quad \times (z_{k_1;\tilde{\alpha}_4} - y_{k_1;\tilde{\alpha}_4}) (z_{k_2;\tilde{\alpha}_5} - y_{k_2;\tilde{\alpha}_5}) (z_{k_3;\tilde{\alpha}_6} - y_{k_3;\tilde{\alpha}_6}) + O\left(\frac{1}{n^2}\right). \end{aligned} \tag{\infty.70}$$

Here, we've again used the standard MSE loss, for which the error factor is given by the residual training error $\epsilon_{j;\tilde{\alpha}} = z_{j;\tilde{\alpha}} - y_{j;\tilde{\alpha}}$. Now, we'll need to pick our generalized Newton tensor $\kappa_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2}$ judiciously in order to make a first update that fully accounts for instantiation-to-instantiation fluctuations of particular networks in our ensemble.

Taking inspiration from our $1/n$ -failed infinite-width Newton's step (∞.62), let's take a similar-looking first step according to

$$\begin{aligned} \eta \kappa_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2} &= \left(\widehat{H}^{-1} \right)_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2} \\ &= \delta_{ij} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \widehat{\Delta H}_{ij;\tilde{\alpha}_3 \tilde{\alpha}_4} \tilde{H}^{\tilde{\alpha}_4 \tilde{\alpha}_2} \\ &\quad + \sum_{k=1}^{n_L} \sum_{\tilde{\alpha}_3, \dots, \tilde{\alpha}_6 \in \mathcal{A}} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \widehat{\Delta H}_{ik;\tilde{\alpha}_3 \tilde{\alpha}_4} \tilde{H}^{\tilde{\alpha}_4 \tilde{\alpha}_5} \widehat{\Delta H}_{kj;\tilde{\alpha}_5 \tilde{\alpha}_6} \tilde{H}^{\tilde{\alpha}_6 \tilde{\alpha}_2} + O(\Delta^3). \end{aligned} \tag{\infty.71}$$

⁹Note that when we developed our interpretation of an infinite-width network as a linear model in §10.4, we made a similar distinction between parameter indices μ and output-component indices i in (10.137) when defining the random feature functions $\hat{\phi}_{i,\mu}(x)$. We also discussed how *wiring* can be incorporated into a nonminimal model of representation learning in §11.4.3.

Here, we've introduced the complete inverse of the stochastic NTK sub-tensor evaluated on the training set, satisfying

$$\sum_{j, \tilde{\alpha}_2} \left(\hat{H}^{-1} \right)_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2} \hat{H}_{jk; \tilde{\alpha}_2 \tilde{\alpha}_3} = \delta_{ik} \delta_{\tilde{\alpha}_3}^{\tilde{\alpha}_1}, \quad (\infty.72)$$

and in the last equality of (∞.71), we've used the Schwinger–Dyson equations (4.55), which is a physicist's way of saying that we expanded the inverse of the stochastic NTK around the NTK mean.¹⁰ The main difference between this new giant leap (∞.71) and our previous infinite-width giant leap (∞.62) is that we're now taking into account the instantiation-to-instantiation fluctuations of the NTK across different realizations of the model parameters; in other words, we're implementing a *different* Newton step for each *particular* network with its associated NTK $\hat{H}_{i_1 i_2; \tilde{\alpha}_1 \tilde{\alpha}_2}$. Accordingly, this step can be thought of as loosely corresponding to the first phase of training for such a particular network.

Taking this step, i.e., plugging this generalized Newton tensor (∞.71) into our update to the network output (∞.70), we find the training error decreases to

$$\begin{aligned} & z_{i; \tilde{\alpha}}(t = 1) - y_{i; \tilde{\alpha}} \\ &= \frac{1}{2} \sum_{\substack{j_1, j_2, \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_4}} \widehat{\mathrm{d}H}_{ij_1 j_2; \tilde{\alpha} \tilde{\alpha}_1 \tilde{\alpha}_2} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_4} (z_{j_1; \tilde{\alpha}_3} - y_{j_1; \tilde{\alpha}_3}) (z_{j_2; \tilde{\alpha}_4} - y_{j_2; \tilde{\alpha}_4}) \\ &\quad - \frac{1}{6} \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_6}} \widehat{\mathrm{dd}_I H}_{ij_1 j_2 j_3; \delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_4} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_6} \\ &\quad \times (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) + O\left(\frac{1}{n^2}\right). \end{aligned} \quad (\infty.73)$$

Thus, we've correctly taken care of the first type of the $1/n$ corrections, and with this step we've reduced the residual prediction error on the training set from the order-one error of our initial prediction,

$$z_{i; \tilde{\alpha}}(t = 0) - y_{i; \tilde{\alpha}} = O(1), \quad (\infty.74)$$

to a much smaller error of

$$z_{i; \tilde{\alpha}}(t = 1) - y_{i; \tilde{\alpha}} = O\left(\frac{1}{n}\right), \quad (\infty.75)$$

loosely corresponding to the empirically-large initial decrease of error when training networks in practice. Moreover, the rest of the error in (∞.73) is now entirely due to the

¹⁰Despite saying that we wouldn't belabor this any further, this is one of those unfortunate situations where we had to decide between decorating the inverse $(\hat{H}^{-1})_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2}$ with either a hat or a tilde, and we went with the hat. Hopefully the tildes on the sample indices, e.g., $\tilde{\alpha}_1, \tilde{\alpha}_2$, will remind you that the inverse of this stochastic object is taken only with respect to the training set. Note that at no point will we ever need the inverse of the stochastic NTK evaluated on a general dataset \mathcal{D} .

additional finite-width corrections encoded by the dNTK and first ddNTK, a consequence of representation learning.¹¹

Second Update: Representation Learning Strikes Back

To reduce the training error further, we'll need to update our network again to further account for the fact that its representation evolved with the first update. In other words, we'll need to make a second gradient-descent update. If you'd like, you can imagine that this update corresponds to a second phase of training – following a first phase where the network significantly decreased its training error with the NTK evaluated at initialization – and so now the model must refine its features in order to further improve its performance.

Accordingly, since our training error is already down to $\sim 1/n$, our second update is going to be a lot smaller than our first. In particular, the update itself will necessarily only be of order $1/n$ so as to precisely cancel the remaining $1/n$ training error; that is, it's actually more of a *small step* than another *giant leap*.¹²

To determine which step we need to take, let's write the network output after a second update as

$$\begin{aligned} z_{i;\delta}(t=2) & \quad (\infty.76) \\ &= z_{i;\delta}(t=1) - \sum_{j,k,\tilde{\alpha}_1,\tilde{\alpha}_2} H_{ij;\delta\tilde{\alpha}_1}(t=1) \eta \kappa_{jk}^{\tilde{\alpha}_1\tilde{\alpha}_2}(t=1) [z_{k;\tilde{\alpha}_2}(t=1) - y_{k;\tilde{\alpha}_2}] + O\left(\frac{1}{n^2}\right) \\ &= z_{i;\delta}(t=1) - \sum_{j,\tilde{\alpha}_1,\tilde{\alpha}_2} H_{\delta\tilde{\alpha}_1} \eta \kappa_{ij}^{\tilde{\alpha}_1\tilde{\alpha}_2}(t=1) [z_{j;\tilde{\alpha}_2}(t=1) - y_{j;\tilde{\alpha}_2}] + O\left(\frac{1}{n^2}\right), \end{aligned}$$

where we left the second update's product of global learning rate and generalized Newton tensor $\eta \kappa_{ij}^{\tilde{\alpha}_1\tilde{\alpha}_2}(t=1)$ unspecified for now. Note here that in the first equality we dropped the d(d)NTK terms from the update: since after our first update ($\infty.73$) the training error has already decreased to $O(1/n)$, the would-be dNTK term is very subleading $\sim \widehat{dH} \times [\epsilon(t=1)]^2 = O(1/n^3)$, and the would-be ddNTK term is ridiculously subleading $\sim \widehat{dd_1H} \times [\epsilon(t=1)]^3 = O(1/n^4)$. Similarly, on the third line we replaced the NTK after the first step $H_{ij;\delta\tilde{\alpha}_1}(t=1)$ by the NTK mean at initialization $\delta_{ij}H_{\delta\tilde{\alpha}_1}$. Given the $1/n$ -suppressed training error after the first step, this substitution can be justified for these two reasons in conjunction: (i) the update to the NTK $\widehat{H}_{ij;\delta\tilde{\alpha}}(t=1) - \widehat{H}_{ij;\delta\tilde{\alpha}}(t=0)$ is itself suppressed by $1/n$, cf. ($\infty.6$), so we may use the version from before the update,

¹¹In particular, this first update ($\infty.71$) would satisfy the training condition ($\infty.66$) only if the NTK were constant under gradient descent. There's actually a name given to this type of phenomenon, *lazy training*, referring to situations when the network function behaves as if it is equal to a linearization around the initial value of its parameters [70]. As we know from our discussion in §10.4, if the NTK is constant, then the network is a linear model.

¹²As we will see, the overall learning rate is essentially the same for both updates; the second update is only smaller because the gradient of the loss after the first update is itself much smaller when near a minimum.

and (ii) the NTK fluctuation is also suppressed compared to its mean, so we may then swap the stochastic NTK at initialization for its mean. This means that if we make the following choice for our *small-step* second update,

$$\eta \kappa_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2}(t=1) = \delta_{ij} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_2} + O\left(\frac{1}{n}\right), \quad (\infty.77)$$

then it's easy to see from (∞.76) that our network will now be fully trained as (∞.66):

$$z_{i;\tilde{\alpha}}(t=2) - y_{i;\tilde{\alpha}} = O\left(\frac{1}{n^2}\right). \quad (\infty.78)$$

Thus, with our second update we were able to reduce the residual training error by an additional factor of $1/n$ as compared to the error after the first update (∞.73).¹³

For general inputs $\delta \in \mathcal{D}$, plugging our choice of learning rate and Newton tensor (∞.77) back into our second update (∞.76) and further re-expressing the network output after the first step $z_{i;\delta}(t=1)$ by (∞.70), with (∞.71) as our choice for the first step, we get

$$\begin{aligned} z_{i;\delta}(t=2) &= z_{i;\delta} - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} H_{\delta \tilde{\alpha}_1} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_2} (z_{i;\tilde{\alpha}_2} - y_{i;\tilde{\alpha}_2}) \\ &\quad + \sum_{j=1}^{n_L} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} \left[\widehat{\Delta H}_{ij;\delta \tilde{\alpha}_1} - \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} H_{\delta \tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_4} \widehat{\Delta H}_{ij;\tilde{\alpha}_4 \tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_2} (z_{j;\tilde{\alpha}_2} - y_{j;\tilde{\alpha}_2}) \\ &\quad - \sum_{j,k=1}^{n_L} \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4 \in \mathcal{A}} \left[\widehat{\Delta H}_{ij;\delta \tilde{\alpha}_1} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6 \in \mathcal{A}} H_{\delta \tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5 \tilde{\alpha}_6} \widehat{\Delta H}_{ij;\tilde{\alpha}_6 \tilde{\alpha}_1} \right] \\ &\quad \quad \times \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_2} \widehat{\Delta H}_{jk;\tilde{\alpha}_2 \tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_4} (z_{k;\tilde{\alpha}_4} - y_{k;\tilde{\alpha}_4}) \\ &\quad + \frac{1}{2} \sum_{j_1, j_2=1}^{n_L} \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4 \in \mathcal{A}} \left[\widehat{\mathrm{d}H}_{ij_1 j_2; \delta \tilde{\alpha}_1 \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6 \in \mathcal{A}} H_{\delta \tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5 \tilde{\alpha}_6} \widehat{\mathrm{d}H}_{ij_1 j_2; \tilde{\alpha}_6 \tilde{\alpha}_1 \tilde{\alpha}_2} \right] \\ &\quad \quad \times \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_4} (z_{j_1; \tilde{\alpha}_3} - y_{j_1; \tilde{\alpha}_3}) (z_{j_2; \tilde{\alpha}_4} - y_{j_2; \tilde{\alpha}_4}) \\ &\quad - \frac{1}{6} \sum_{j_1, j_2, j_3=1}^{n_L} \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_6 \in \mathcal{A}} \left[\widehat{\mathrm{dd}_1 H}_{ij_1 j_2 j_3; \delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8 \in \mathcal{A}} H_{\delta \tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7 \tilde{\alpha}_8} \widehat{\mathrm{dd}_1 H}_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] \\ &\quad \quad \times \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_4} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_6} (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) \\ &\quad + O\left(\frac{1}{n^2}\right). \end{aligned} \quad (\infty.79)$$

¹³This suggests that additional updates could continue to reduce the training error by additional factors of $1/n$. However, these further refinements – determined in terms of the higher-order corrections to our effective theory description – will be qualitatively the same as our leading finite-width description at $O(L/n)$. In other words, improving an infinite-width description to finite-width description incorporates representation learning, while more precise finite-width descriptions just allow the model to make further refinements to its features. Practically speaking, we expect our leading finite-width description to be very accurate for networks with reasonable values of the aspect ratio L/n .

Here, we see that the expressions in the square brackets vanish identically for all the training inputs $\delta = \tilde{\alpha} \in \mathcal{A}$: our network is thus fully trained. In particular, since all of the variables on the right-hand side of (∞.79) are at initialization, this solution realizes our goal (∞.55) of expressing the output of a fully-trained network as a functional of such variables at initialization. Accordingly, the statistics of the fully-trained distribution (∞.56) can now be worked out from the joint preactivation–NTK–dNTK–ddNTKs distribution (∞.53) at initialization.¹⁴

While this is all very exciting, let us also caution you that these generalized second-order updates are probably best thought of as giving a simple theoretical model of a training algorithm – designed to let us understand the training process analytically – and by no means are we suggesting that they provide a good or useful option for practical optimization. Instead, the most practical algorithm for optimization is vanilla first-order gradient descent. As we’ve already pointed out that the output of a fully-trained network *does* depend on the details of the algorithm used for optimization, now we really have no choice left other than to explicitly analyze many many steps of gradient descent.

¹⁴Alternatively, we could have reached this same solution in a *single update* if we had instead made a *finely-tuned* giant leap, picking the generalized Newton tensor as

$$\begin{aligned} \eta \kappa_{ij}^{\tilde{\alpha}_1 \tilde{\alpha}_2} &= \delta_{ij} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \widehat{\Delta H}_{ij; \tilde{\alpha}_3 \tilde{\alpha}_4} \tilde{H}^{\tilde{\alpha}_4 \tilde{\alpha}_2} \\ &+ \sum_{k=1}^{n_L} \sum_{\tilde{\alpha}_3, \dots, \tilde{\alpha}_6 \in \mathcal{A}} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \widehat{\Delta H}_{ik; \tilde{\alpha}_3 \tilde{\alpha}_4} \tilde{H}^{\tilde{\alpha}_4 \tilde{\alpha}_5} \widehat{\Delta H}_{kj; \tilde{\alpha}_5 \tilde{\alpha}_6} \tilde{H}^{\tilde{\alpha}_6 \tilde{\alpha}_2} \\ &+ \frac{1}{2} \sum_{k=1}^{n_L} \sum_{\tilde{\alpha}_3, \dots, \tilde{\alpha}_6 \in \mathcal{A}} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_4} \tilde{H}^{\tilde{\alpha}_5 \tilde{\alpha}_6} \widehat{\mathrm{d}H}_{ijk; \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5} (z_{k; \tilde{\alpha}_6} - y_{k; \tilde{\alpha}_6}) \\ &- \frac{1}{6} \sum_{k_1, k_2=1}^{n_L} \sum_{\tilde{\alpha}_3, \dots, \tilde{\alpha}_8 \in \mathcal{A}} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_4} \tilde{H}^{\tilde{\alpha}_5 \tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_6 \tilde{\alpha}_8} \widehat{\mathrm{dd}_1 H}_{ijk_1 k_2; \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \\ &\quad \times (z_{k_1; \tilde{\alpha}_7} - y_{k_1; \tilde{\alpha}_7}) (z_{k_2; \tilde{\alpha}_8} - y_{k_2; \tilde{\alpha}_8}). \end{aligned} \quad (\infty.80)$$

In essence, the first three terms of (∞.80) come from inverting the stochastic NTK, corresponding to our first-update giant leap (∞.71) and account for the instantiation-to-instantiation fluctuations in the NTK. In contrast, the final two terms correspond to our second-update small step (∞.78) and account for the dNTK–ddNTK-induced representation learning. Note that this generalized Newton tensor (∞.80) is *asymmetric* under the exchange of paired indices $(i, \tilde{\alpha}_1) \leftrightarrow (j, \tilde{\alpha}_2)$ due to the last term, and thus this finely-tuned update doesn’t admit an alternative interpretation of optimization with a further generalized loss (∞.69).

This single update is analogous to the *direct optimization* solution (11.132) for quadratic regression in §11.4. The very finely-tuned nature of this single-update algorithm (∞.80) suggests that it’s easier to make the fine adjustments required to reach a solution by taking many simpler steps rather than fewer complicated steps. We will see the other side of this next in §∞.2.2 when we study training by many many steps of vanilla gradient descent.

∞ .2.2 Many Many Steps of Gradient Descent

In this extended subsection, we're going to study tensorial gradient descent (7.11) and optimize finite-width neural networks according to the MSE loss with a constant global learning rate η .¹⁵ With this progenitor-of-all-other-gradient-based-learning-algorithms algorithm, the network output will evolve as (∞ .4),

$$\begin{aligned} z_{i;\delta}(t+1) &= z_{i;\delta}(t) - \eta \sum_{j,\tilde{\alpha}} H_{ij;\delta\tilde{\alpha}}(t) [z_{j;\tilde{\alpha}}(t) - y_{j;\tilde{\alpha}}] \\ &\quad + \frac{\eta^2}{2} \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2} dH_{ij_1j_2;\delta\tilde{\alpha}_1\tilde{\alpha}_2}(t) [z_{j_1;\tilde{\alpha}_1}(t) - y_{j_1;\tilde{\alpha}_1}] [z_{j_2;\tilde{\alpha}_2}(t) - y_{j_2;\tilde{\alpha}_2}] \\ &\quad - \frac{\eta^3}{6} \sum_{j_1,j_2,j_3,\tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3} \widehat{\text{dd}_I H}_{ij_1j_2j_3;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} \\ &\quad \times [z_{j_1;\tilde{\alpha}_1}(t) - y_{j_1;\tilde{\alpha}_1}] [z_{j_2;\tilde{\alpha}_2}(t) - y_{j_2;\tilde{\alpha}_2}] [z_{j_3;\tilde{\alpha}_3}(t) - y_{j_3;\tilde{\alpha}_3}] ; \end{aligned} \quad (\infty.81)$$

in conjunction the NTK will evolve as (∞ .6),

$$\begin{aligned} H_{i_1i_2;\delta_1\delta_2}(t+1) &= H_{i_1i_2;\delta_1\delta_2}(t) - \eta \sum_{j,\tilde{\alpha}} \left(dH_{i_1i_2j;\delta_1\delta_2\tilde{\alpha}}(t) + dH_{i_2i_1j;\delta_2\delta_1\tilde{\alpha}}(t) \right) [z_{j;\tilde{\alpha}}(t) - y_{j;\tilde{\alpha}}] \\ &\quad + \frac{\eta^2}{2} \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2} \left(\widehat{\text{dd}_I H}_{i_1i_2j_1j_2;\delta_1\delta_2\tilde{\alpha}_1\tilde{\alpha}_2} + \widehat{\text{dd}_I H}_{i_2i_1j_1j_2;\delta_2\delta_1\tilde{\alpha}_1\tilde{\alpha}_2} + 2\widehat{\text{dd}_{II} H}_{i_1i_2j_1j_2;\delta_1\delta_2\tilde{\alpha}_1\tilde{\alpha}_2} \right) \\ &\quad \times [z_{j_1;\tilde{\alpha}_1}(t) - y_{j_1;\tilde{\alpha}_1}] [z_{j_2;\tilde{\alpha}_2}(t) - y_{j_2;\tilde{\alpha}_2}] ; \end{aligned} \quad (\infty.82)$$

and in further conjunction the dNTK will evolve as (∞ .9),

$$\begin{aligned} dH_{i_0i_1i_2;\delta_0\delta_1\delta_2}(t+1) &= dH_{i_0i_1i_2;\delta_0\delta_1\delta_2}(t) \\ &\quad - \eta \sum_{j,\tilde{\alpha}} \left(\widehat{\text{dd}_I H}_{i_0i_1i_2j;\delta_0\delta_1\delta_2\tilde{\alpha}} + \widehat{\text{dd}_{II} H}_{i_0i_1i_2j;\delta_0\delta_1\delta_2\tilde{\alpha}} + \widehat{\text{dd}_{II} H}_{i_0i_2i_1j;\delta_0\delta_2\delta_1\tilde{\alpha}} \right) [z_{j;\tilde{\alpha}}(t) - y_{j;\tilde{\alpha}}] . \end{aligned} \quad (\infty.83)$$

In writing these dynamical equations, we've stopped explicitly denoting that our equations have $O(1/n^2)$ errors, and we've also expressed the fact that the ddNTKs are

¹⁵Don't let the word *tensorial* scare you here; this just means that we will allow for our training hyperparameters $\lambda_b^{(\ell)}$ and $\lambda_W^{(\ell)}$ as part of the definition of the NTK. We hope it's already clear why including these hyperparameters is a good idea – if not, please flip back to §9.4 and reread the paragraphs on the learning rate equivalence principle – and they are actually simple to include practically as part of any optimization algorithm.

Moreover, as they are just part of the definition of the NTK, they have absolutely no consequence on the dynamics presented here; i.e., our solution also covers *nontensorial* gradient descent, though in such a case you'd have different asymptotic solutions for the statistics of the NTK, dNTK, and ddNTKs. This is also why we think of the training hyperparameters $\lambda_b^{(\ell)}$ and $\lambda_W^{(\ell)}$ as being independent from the details of the optimization algorithm itself.

t -independent at order $1/n$, using their values at initialization, $\widehat{\text{dd}}_{\text{I}} H_{i_1 i_2 i_3 i_4; \delta_1 \delta_2 \delta_3 \delta_4}$ and $\widehat{\text{dd}}_{\text{II}} H_{i_1 i_2 i_3 i_4; \delta_1 \delta_2 \delta_3 \delta_4}$, with hats to remind us of their stochasticity. These joint updates (∞.81), (∞.82), and (∞.83) are coupled *difference equations*, and the equations for the network output and the dynamical NTK are both *nonlinear* in the output $z_{i;\delta}$.

Although this seems daunting, we are now going to solve these equations in a closed form. First, we'll analyze the dynamics while neglecting the finite-width effects of the dNTK and ddNTKs, in which the problem will reduce to a single *linear* difference equation. Then, we'll use perturbation theory to incorporate the effect of all the NTK differentials and allow for a dynamically evolving NTK and dNTK.

Free Theory: Step-Independent NTK

Let's begin by setting the dNTK and ddNTKs to zero. Since their leading statistics are $1/n$ -suppressed, we'll be able to use perturbation theory to reincorporate all their effects later. In this limit the NTK update equation (∞.82) is trivial, solved by the *free* or *step-independent* NTK:

$$H_{i_1 i_1; \delta_1 \delta_2}(t) = H_{i_1 i_2; \delta_1 \delta_2}(t=0) \equiv \hat{H}_{i_1 i_2; \delta_1 \delta_2}. \quad (\infty.84)$$

Unsurprisingly, this just means that when the dNTK and ddNTKs vanish, the NTK doesn't update from its initialization.

Plugging this solution into the preactivation update equation (∞.81) and turning off the dNTK and ddNTKs, the remaining dynamical equation simplifies to

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) - \eta \sum_{j,\tilde{\alpha}} \hat{H}_{ij;\delta\tilde{\alpha}} [z_{j;\tilde{\alpha}}(t) - y_{j;\tilde{\alpha}}]. \quad (\infty.85)$$

Thus, the residual training error, $z_{\tilde{\alpha}}(t) - y_{\tilde{\alpha}}$, sources the updates to the network output $z_{\delta}(t)$ for general inputs $\delta \in \mathcal{D}$. Moreover, when restricted to the inputs from the training set $\tilde{\alpha} \in \mathcal{A}$, we can rewrite this difference equation (∞.85) as

$$z_{i;\tilde{\alpha}}(t+1) - y_{i;\tilde{\alpha}} = \sum_{j,\tilde{\alpha}_1} \left(I_{ij;\tilde{\alpha}\tilde{\alpha}_1} - \eta \hat{H}_{ij;\tilde{\alpha}\tilde{\alpha}_1} \right) [z_{j;\tilde{\alpha}_1}(t) - y_{j;\tilde{\alpha}_1}], \quad (\infty.86)$$

where we've defined the *identity operator*

$$I_{i_1 i_2; \tilde{\alpha}_1 \tilde{\alpha}_2} \equiv \delta_{i_1 i_2} \delta_{\tilde{\alpha}_1 \tilde{\alpha}_2}. \quad (\infty.87)$$

In this form, (∞.86) is a first-order homogeneous linear difference equation for the residual training error, $z_{\tilde{\alpha}}(t) - y_{\tilde{\alpha}}$, which is just a fancy way of saying that this is going to be a piece of cake.

In particular, the update to the prediction error is just a simple multiplication by a constant matrix, and the solution is given by an exponential:

$$z_{i;\tilde{\alpha}}^{\text{F}}(t) - y_{i;\tilde{\alpha}} = \sum_{j,\tilde{\alpha}_1} U_{ij;\tilde{\alpha}\tilde{\alpha}_1}(t) (z_{j;\tilde{\alpha}_1} - y_{j;\tilde{\alpha}_1}). \quad (\infty.88)$$

Here, on the left-hand side, we've labeled the solution with an "F" to indicate it's the *free solution*, with the nonlinear effects from the dNTK and ddNTKs turned off; on the right-hand side, we have the residual training error at initialization, and the **step-evolution operator** is defined as an iterative product of t steps:

$$\begin{aligned} U_{i_0 i_0; \tilde{\alpha}_0 \tilde{\alpha}_0}(t) &\equiv \left[(I - \eta \hat{H})^t \right]_{i_0 i_0; \tilde{\alpha}_0 \tilde{\alpha}_0} \\ &= \sum_{\substack{i_1, \dots, i_{t-1} \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_{t-1}}} \left(I_{i_t i_{t-1}; \tilde{\alpha}_t \tilde{\alpha}_{t-1}} - \eta \hat{H}_{i_t i_{t-1}; \tilde{\alpha}_t \tilde{\alpha}_{t-1}} \right) \cdots \left(I_{i_1 i_0; \tilde{\alpha}_1 \tilde{\alpha}_0} - \eta \hat{H}_{i_1 i_0; \tilde{\alpha}_1 \tilde{\alpha}_0} \right). \end{aligned} \quad (\infty.89)$$

For any positive-definite NTK and sufficiently small global learning rate η , this operator will exponentially decay to zero, $U(t) \rightarrow 0$, as the number of steps becomes large, $t \rightarrow \infty$.¹⁶ Thus, the residual training error ($\infty.88$) will vanish exponentially quickly,

$$\lim_{t \rightarrow \infty} z_{i; \tilde{\alpha}}^F(t) = y_{i; \tilde{\alpha}}, \quad (\infty.90)$$

with the step scale for the decay of the individual components set by the step-independent NTK.¹⁷

Having now solved the free dynamics on the training set, we can plug this solution ($\infty.88$) back into the difference equation ($\infty.85$) for general inputs $\delta \in \mathcal{D}$. With the source known explicitly, we can easily write down a solution that satisfies the initial condition $z_{i; \delta}^F(t=0) = z_{i; \delta}$:

$$\begin{aligned} z_{i; \delta}^F(t) &= z_{i; \delta} - \sum_{j, \tilde{\alpha}} \hat{H}_{ij; \delta \tilde{\alpha}} \left\{ \eta \sum_{s=0}^{t-1} \left[z_{j; \tilde{\alpha}}^F(s) - y_{j; \tilde{\alpha}} \right] \right\} \\ &= z_{i; \delta} - \sum_{j, \tilde{\alpha}_1, \tilde{\alpha}_2} \hat{H}_{ij; \delta \tilde{\alpha}_1} a_{j; \tilde{\alpha}_2}(t). \end{aligned} \quad (\infty.91)$$

¹⁶In particular, for this limit to converge we just need $\|I - \eta \hat{H}\|_\infty < 1$, i.e., the largest eigenvalue of the operator $I - \eta \hat{H}$ must be less than one. With our attention to the principles of criticality and equivalence, our choices of initialization and training hyperparameters were made so that the NTK is always of order one, and thus it's very easy for our networks to satisfy this constraint.

¹⁷If we wanted to study the ODE or *continuum limit* of the dynamics, we could take the global learning rate to zero, $\eta \rightarrow 0$, while holding the product, $\tau \equiv \eta t$, fixed. In such a limit, the step-evolution operator becomes simply $U(t) \rightarrow \exp(-\hat{H}\tau)$. While such a limit is mostly unnecessary for any theoretical purpose – it's just as easy to study the discrete dynamics that actually describe the practical optimization algorithm – it does provide more substance to the objection in footnote 5 of §7.2 of the name *neural tangent kernel* for the stochastic operator \hat{H} . In particular, this continuum limit makes clear that the NTK is really best thought of as a *Hamiltonian* as it generates the evolution of observables, and the step-evolution operator $U(t)$ is like a unitary time-evolution operator, albeit in *imaginary* time. More precisely, in the limit where the dNTK and ddNTKs are set to zero, the NTK is akin to a *free* Hamiltonian, with exactly solvable dynamics; with a nonzero dNTK or ddNTKs, the Hamiltonian includes nontrivial *interactions* and can be analyzed via time-dependent perturbation theory.

Here we've defined a dynamical helper function, with an explicit representation given by

$$\begin{aligned}
 a_{j;\tilde{\alpha}}(t) &\equiv \eta \sum_{s=0}^{t-1} \left[z_{j;\tilde{\alpha}}^{\text{F}}(s) - y_{j;\tilde{\alpha}} \right] = \eta \sum_{s=0}^{t-1} \left[\sum_{k;\tilde{\alpha}_1} U_{jk;\tilde{\alpha}\tilde{\alpha}_1}(t) (z_{k;\tilde{\alpha}_1} - y_{k;\tilde{\alpha}_1}) \right] \\
 &= \eta \sum_{k;\tilde{\alpha}_1} \left\{ \sum_{s=0}^{t-1} \left[(I - \eta \hat{H})^s \right]_{jk;\tilde{\alpha}\tilde{\alpha}_1} \right\} (z_{k;\tilde{\alpha}_1} - y_{k;\tilde{\alpha}_1}) \\
 &= \eta \sum_{k,m;\tilde{\alpha}_1,\tilde{\alpha}_2} \left\{ \left[I - (I - \eta \hat{H}) \right]^{-1} \right\}_{jm}^{\tilde{\alpha}\tilde{\alpha}_2} \left[I - (I - \eta \hat{H}) \right]_{mk;\tilde{\alpha}_2\tilde{\alpha}_1}^t (z_{k;\tilde{\alpha}_1} - y_{k;\tilde{\alpha}_1}) \\
 &= \sum_{m,\tilde{\alpha}_2} (\hat{H}^{-1})_{jm}^{\tilde{\alpha}\tilde{\alpha}_2} \left\{ z_{m;\tilde{\alpha}_2} - y_{m;\tilde{\alpha}_2} - \left[z_{m;\tilde{\alpha}_2}^{\text{F}}(t) - y_{m;\tilde{\alpha}_2} \right] \right\}.
 \end{aligned} \tag{\infty.92}$$

In this expression, to get to the third line we used the standard formula for evaluating geometric sums, $1 + x + x^2 + \dots + x^{t-1} = (1 - x^t)/(1 - x)$, and to get to the final line we evaluated the inverse and substituted in for the step-evolution operator (∞.89). Recall also that \hat{H}^{-1} , defined by (∞.72), is the inverse of the stochastic NTK submatrix evaluated on the training set for a particular realization of the parameters.¹⁸

In particular, for sufficiently small η , as the number of steps becomes large, $t \rightarrow \infty$, the residual error $z_{m;\tilde{\alpha}_2}^{\text{F}}(t) - y_{m;\tilde{\alpha}_2}$ exponentially vanishes (∞.90). Thus, the prediction for a general input $\delta \in \mathcal{D}$ (∞.91) will exponentially converge to

$$z_{i;\delta}^{\text{F}}(t = \infty) = z_{i;\delta} - \sum_{j,k;\tilde{\alpha}_1,\tilde{\alpha}_2} \hat{H}_{ij;\delta\tilde{\alpha}_1} (\hat{H}^{-1})_{jk}^{\tilde{\alpha}_1\tilde{\alpha}_2} (z_{k;\tilde{\alpha}_2} - y_{k;\tilde{\alpha}_2}). \tag{\infty.93}$$

Although we took the limit of a large number of steps here, the exponential convergence means that for any number of steps T such that $T \gtrsim (\eta \hat{H})^{-1}$, the prediction error will be exponentially close to its final $T \rightarrow \infty$ value.

In fact, this solution (∞.93) precisely matches the solution we would have gotten in §∞.2.1 after the first update (∞.68) with the Newton tensor (∞.71), so long as the dNTK and ddNTKs are turned off. This means that the *algorithm dependence* that emerges at finite width is solely due to the presence of the NTK differentials and the resulting nonlinear dynamics.

Interacting Theory: Dynamical NTK and dNTK

Now, let's incorporate the nonzero dNTK and ddNTKs into our analysis. To do so, we need to decompose both the dynamical network output and the dynamical NTK as

$$z_{i;\delta}(t) \equiv z_{i;\delta}^{\text{F}}(t) + z_{i;\delta}^{\text{I}}(t), \tag{\infty.94}$$

$$H_{ij;\delta\tilde{\alpha}}(t) \equiv \hat{H}_{ij;\delta\tilde{\alpha}} + H_{ij;\delta\tilde{\alpha}}^{\text{I}}(t). \tag{\infty.95}$$

¹⁸Unfortunately, in the absence of a Newton tensor, the raised/lowered sample indices in (∞.92) do not align well. (In fact, the problem really began in the update equation (∞.85) with the doubled-lowered $\tilde{\alpha}$ index.) If you'd like, you can fix this by judicious use of identity matrices such as $\delta_{\tilde{\alpha}_1\tilde{\alpha}_2}$ and $\delta^{\tilde{\alpha}_3\tilde{\alpha}_4}$. However, such usage over-clutters the presentation and so is probably not worth it, despite the additional clarity and *type safety* that such index alignment provides.

Here, for the network output, $z_{i;\delta}^F(t)$ is the *free* part, satisfying the linear difference equation (∞.85) with a solution given by (∞.91), and $z_{i;\delta}^I(t)$ is the *interacting* part, encapsulating the corrections due to all the nonzero NTK differentials. Similarly, for the NTK, $\widehat{H}_{ij;\delta\tilde{\alpha}}$ is the step-independent or *free* part, fixed at initialization, and $H_{ij;\delta\tilde{\alpha}}^I(t)$ is the dynamical step-dependent or *interaction* NTK. Since both $z_{i;\delta}^I(t)$ and $H_{ij;\delta\tilde{\alpha}}^I(t)$ are absent in the free limit, $\widehat{dH}, \widehat{dd_I H}, \widehat{dd_{II} H} \rightarrow 0$, at leading order we expect them to be a linear combination of these objects, schematically

$$z^I(t) = [\text{thing } 0](t) \widehat{dH} + [\text{thing } 1](t) \widehat{dd_I H} + [\text{thing } 2](t) \widehat{dd_{II} H}, \quad (\infty.96)$$

$$H^I(t) = [\text{thing } 0](t) \widehat{dH} + [\text{thing } 1](t) \widehat{dd_I H} + [\text{thing } 2](t) \widehat{dd_{II} H}, \quad (\infty.97)$$

where various tensorial things will have various time dependencies. This means in turn that any product of $z_{i;\delta}^I(t)$ or $H_{ij;\delta\tilde{\alpha}}^I(t)$ with one of $\widehat{dH}, \widehat{dd_I H}, \widehat{dd_{II} H}$ can be neglected to leading order, which is the main reason why we'll be able to systematically solve these nonlinear dynamics by perturbation theory. Finally, the initial condition for the interacting parts of the network output and the NTK must satisfy

$$z_{i;\delta}^I(t=0) = 0, \quad (\infty.98)$$

$$H_{ij;\delta\tilde{\alpha}}^I(t=0) = 0, \quad (\infty.99)$$

since the free part of the network output already satisfies $z_{i;\delta}^F(t=0) = z_{i;\delta}$ at initialization, and the free part of the NTK is step-independent and thus *is* the NTK at initialization. With all those in mind, our goal now is to find a solution for $z_{i;\delta}^I(t)$ such that the full network output, $z_{i;\delta}(t)$, the interaction NTK, $H_{ij;\delta\tilde{\alpha}}^I(t)$, and the dynamical dNTK, $dH_{i_0 i_1 i_2; \delta_0 \delta_1 \delta_2}(t)$, all together satisfy the coupled nonlinear dynamics (∞.81), (∞.82), and (∞.83) to leading order in $1/n$.

First, let's work out the dynamics of the dNTK. From (∞.83) we see that the updates to the dNTK are sourced by the residual training error $z_{i;\tilde{\alpha}}(t) - y_{j;\tilde{\alpha}}$. Using our decomposition for the network output, (∞.94), and the initial condition,

$$dH_{i_0 i_1 i_2; \delta_0 \delta_1 \delta_2}(t=0) = \widehat{dH}_{i_0 i_1 i_2; \delta_0 \delta_1 \delta_2}, \quad (\infty.100)$$

after iterating the dynamics, we have

$$\begin{aligned} dH_{i_0 i_1 i_2; \delta_0 \delta_1 \delta_2}(t) &= \widehat{dH}_{i_0 i_1 i_2; \delta_0 \delta_1 \delta_2} - \sum_{j, \tilde{\alpha}} \left(\widehat{dd_I H}_{i_0 i_1 i_2 j; \delta_0 \delta_1 \delta_2 \tilde{\alpha}} + \widehat{dd_{II} H}_{i_0 i_1 i_2 j; \delta_0 \delta_1 \delta_2 \tilde{\alpha}} + \widehat{dd_{III} H}_{i_0 i_2 i_1 j; \delta_0 \delta_2 \delta_1 \tilde{\alpha}} \right) \\ &\quad \times \left\{ \eta \sum_{s=0}^{t-1} \left[z_{j;\tilde{\alpha}}^F(s) + z_{j;\tilde{\alpha}}^I(s) - y_{j;\tilde{\alpha}} \right] \right\} \\ &= \widehat{dH}_{i_0 i_1 i_2; \delta_0 \delta_1 \delta_2} - \sum_{j, \tilde{\alpha}} \left(\widehat{dd_I H}_{i_0 i_1 i_2 j; \delta_0 \delta_1 \delta_2 \tilde{\alpha}} + \widehat{dd_{II} H}_{i_0 i_1 i_2 j; \delta_0 \delta_1 \delta_2 \tilde{\alpha}} + \widehat{dd_{III} H}_{i_0 i_2 i_1 j; \delta_0 \delta_2 \delta_1 \tilde{\alpha}} \right) a_{j;\tilde{\alpha}}(t). \end{aligned} \quad (\infty.101)$$

Here in the last step, we first dropped the product of $z_{i;\delta}^I(t)$ with $\widehat{dd_I H}$, as explained earlier, and then substituted in our dynamical helper function (∞.92). Now, plugging in

our final expression from (∞.92) and neglecting the fluctuation part of the NTK inverse as subleading, we find an expression for the dynamical dNTK:

$$\begin{aligned} \widehat{dH}_{i_0 i_1 i_2; \delta_0 \delta_1 \delta_2}(t) = & \widehat{dH}_{i_0 i_1 i_2; \delta_0 \delta_1 \delta_2} - \sum_{k, \tilde{\alpha}_1, \tilde{\alpha}_2} \left(\widehat{dd_I H}_{i_0 i_1 i_2 j; \delta_0 \delta_1 \delta_2 \tilde{\alpha}_1} \right. \\ & + \widehat{dd_{II} H}_{i_0 i_1 i_2 j; \delta_0 \delta_1 \delta_2 \tilde{\alpha}_1} + \widehat{dd_{II} H}_{i_0 i_2 i_1 j; \delta_0 \delta_2 \delta_1 \tilde{\alpha}_1} \Big) \\ & \times \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_2} \left\{ z_{j; \tilde{\alpha}_2} - y_{j; \tilde{\alpha}_2} - \left[z_{j; \tilde{\alpha}_2}^F(t) - y_{j; \tilde{\alpha}_2} \right] \right\}. \end{aligned} \quad (\infty.102)$$

In particular, the quantity in the curly brackets represents the difference in training errors between initialization and step t : the larger this difference – i.e., the more the residual training error decreases – the greater the evolution of the dNTK and the more the meta feature functions evolve, undergoing their own form of *meta representation learning* over the course of training.

Next, using our decompositions for the network output and NTK, (∞.94) and (∞.95), we can rewrite the NTK dynamics (∞.82) as a difference equation for the interaction NTK:

$$\begin{aligned} H_{i_1 i_2; \delta_1 \delta_2}^I(t+1) & \quad (\infty.103) \\ = & H_{i_1 i_2; \delta_1 \delta_2}^I(t) - \eta \sum_{j, \tilde{\alpha}} \left[dH_{i_1 i_2 j; \delta_1 \delta_2 \tilde{\alpha}}(t) + dH_{i_2 i_1 j; \delta_2 \delta_1 \tilde{\alpha}}(t) \right] \left[z_{j; \tilde{\alpha}}^F(t) - y_{j; \tilde{\alpha}} \right] \\ & + \frac{\eta^2}{2} \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2} \left(\widehat{dd_I H}_{i_1 i_2 j_1 j_2; \delta_1 \delta_2 \tilde{\alpha}_1 \tilde{\alpha}_2} + \widehat{dd_I H}_{i_2 i_1 j_1 j_2; \delta_2 \delta_1 \tilde{\alpha}_1 \tilde{\alpha}_2} + 2\widehat{dd_{II} H}_{i_1 i_2 j_1 j_2; \delta_1 \delta_2 \tilde{\alpha}_1 \tilde{\alpha}_2} \right) \\ & \times \left[z_{j_1; \tilde{\alpha}_1}^F(t) - y_{j_1; \tilde{\alpha}_1} \right] \left[z_{j_2; \tilde{\alpha}_2}^F(t) - y_{j_2; \tilde{\alpha}_2} \right]. \end{aligned}$$

Here, once again, we've dropped the interacting part of the network output on the right-hand side, as it is always multiplied by either the dNTK or the ddNTKs and thus will be subleading in $1/n$. Denoting the free part of the residual training error (∞.88) as

$$\epsilon_{j; \tilde{\alpha}}^F(t) \equiv z_{j; \tilde{\alpha}}^F(t) - y_{j; \tilde{\alpha}} \quad (\infty.104)$$

for notational convenience and substituting in our solution for the dynamical dNTK (∞.102), we get

$$\begin{aligned} H_{i_1 i_2; \delta_1 \delta_2}^I(t+1) - H_{i_1 i_2; \delta_1 \delta_2}^I(t) & \quad (\infty.105) \\ = & -\eta \sum_{j, \tilde{\alpha}} \left(\widehat{dH}_{i_1 i_2 j; \delta_1 \delta_2 \tilde{\alpha}} + \widehat{dH}_{i_2 i_1 j; \delta_2 \delta_1 \tilde{\alpha}} \right) \epsilon_{j; \tilde{\alpha}}^F(t) \\ & + \eta \sum_{j, k, \tilde{\alpha}, \tilde{\alpha}_1, \tilde{\alpha}_2} \left(\widehat{dd_I H}_{i_1 i_2 j k; \delta_1 \delta_2 \tilde{\alpha} \tilde{\alpha}_1} + \widehat{dd_{II} H}_{i_1 i_2 j k; \delta_1 \delta_2 \tilde{\alpha} \tilde{\alpha}_1} + \widehat{dd_{II} H}_{i_1 j i_2 k; \delta_1 \tilde{\alpha} \delta_2 \tilde{\alpha}_1} \right. \\ & \quad + \widehat{dd_I H}_{i_2 i_1 j k; \delta_2 \delta_1 \tilde{\alpha} \tilde{\alpha}_1} + \widehat{dd_{II} H}_{i_2 i_1 j k; \delta_2 \delta_1 \tilde{\alpha} \tilde{\alpha}_1} + \widehat{dd_{II} H}_{i_2 j i_1 k; \delta_2 \tilde{\alpha} \delta_1 \tilde{\alpha}_1} \Big) \\ & \quad \times \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_2} \epsilon_{j; \tilde{\alpha}}^F(t) \left[z_{k; \tilde{\alpha}_2} - y_{k; \tilde{\alpha}_2} - \epsilon_{k; \tilde{\alpha}_2}^F(t) \right] \\ & + \frac{\eta^2}{2} \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2} \left(\widehat{dd_I H}_{i_1 i_2 j_1 j_2; \delta_1 \delta_2 \tilde{\alpha}_1 \tilde{\alpha}_2} + \widehat{dd_I H}_{i_2 i_1 j_1 j_2; \delta_2 \delta_1 \tilde{\alpha}_1 \tilde{\alpha}_2} + 2\widehat{dd_{II} H}_{i_1 i_2 j_1 j_2; \delta_1 \delta_2 \tilde{\alpha}_1 \tilde{\alpha}_2} \right) \\ & \quad \times \epsilon_{j_1; \tilde{\alpha}_1}^F(t) \epsilon_{j_2; \tilde{\alpha}_2}^F(t). \end{aligned}$$

In particular, the step dependence on the right-hand side is expressed entirely in terms of the free residual training error $\epsilon_{j;\tilde{\alpha}}^F(t)$, and each term is either linear or quadratic in $\epsilon_{j;\tilde{\alpha}}^F(t)$. Thus, in order to solve this difference equation and get $H_{i_1 i_2; \delta_1 \delta_2}^I(t)$, we'll just have to compute sums over these terms.

One of those sums – the one that's linear in the free residual training error – is the dynamical helper function $a_{j;\tilde{\alpha}}(t) \equiv \eta \sum_{s=0}^{t-1} \epsilon_{j;\tilde{\alpha}}^F(t)$ that we evaluated in (∞.92). The other type of sum is quadratic in the free residual training error, which will define a second dynamical helper function:

$$\begin{aligned} b_{j_1 j_2; \tilde{\alpha}_1 \tilde{\alpha}_2}(t) &\equiv \eta \sum_{s=0}^{t-1} \epsilon_{j_1; \tilde{\alpha}_1}^F(t) \epsilon_{j_2; \tilde{\alpha}_2}^F(t) \\ &= \eta \sum_{k_1, k_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \sum_{s=0}^{t-1} \left[(I - \eta \hat{H})^s \right]_{j_1 k_1; \tilde{\alpha}_1 \tilde{\alpha}_3} \\ &\quad \left[(I - \eta \hat{H})^s \right]_{j_2 k_2; \tilde{\alpha}_2 \tilde{\alpha}_4} (z_{k_1; \tilde{\alpha}_3} - y_{k_1; \tilde{\alpha}_3}) (z_{k_2; \tilde{\alpha}_4} - y_{k_2; \tilde{\alpha}_4}) \\ &= \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4} X_{\Pi}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \left[(z_{j_1; \tilde{\alpha}_3} - y_{j_1; \tilde{\alpha}_3}) (z_{j_2; \tilde{\alpha}_4} - y_{j_2; \tilde{\alpha}_4}) - \epsilon_{j_1; \tilde{\alpha}_3}^F(t) \epsilon_{j_2; \tilde{\alpha}_4}^F(t) \right] \\ &\quad + O\left(\frac{1}{n}\right). \end{aligned} \tag{\infty.106}$$

To evaluate this sum, we again used our expression for the free residual training error, (∞.88), in conjunction with the definition of the step-evolution operator, (∞.89). Then, we replaced the stochastic NTK $\hat{H}_{ij; \tilde{\alpha}_1 \tilde{\alpha}_2}$ of the training set by its mean $\delta_{ij} \tilde{H}_{\tilde{\alpha}_1 \tilde{\alpha}_2}$, at the cost of subleading corrections, and formally performed the geometric sum as in (∞.92). This last operation yielded an *inverting tensor* $X_{\Pi}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4}$ implicitly defined by

$$\begin{aligned} \delta_{\tilde{\alpha}_5}^{\tilde{\alpha}_1} \delta_{\tilde{\alpha}_6}^{\tilde{\alpha}_2} &= \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} X_{\Pi}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \frac{1}{\eta} \left[\delta_{\tilde{\alpha}_3 \tilde{\alpha}_5} \delta_{\tilde{\alpha}_4 \tilde{\alpha}_6} - (\delta_{\tilde{\alpha}_3 \tilde{\alpha}_5} - \eta \tilde{H}_{\tilde{\alpha}_3 \tilde{\alpha}_5}) (\delta_{\tilde{\alpha}_4 \tilde{\alpha}_6} - \eta \tilde{H}_{\tilde{\alpha}_4 \tilde{\alpha}_6}) \right] \\ &= \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} X_{\Pi}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \left(\tilde{H}_{\tilde{\alpha}_3 \tilde{\alpha}_5} \delta_{\tilde{\alpha}_4 \tilde{\alpha}_6} + \delta_{\tilde{\alpha}_3 \tilde{\alpha}_5} \tilde{H}_{\tilde{\alpha}_4 \tilde{\alpha}_6} - \eta \tilde{H}_{\tilde{\alpha}_3 \tilde{\alpha}_5} \tilde{H}_{\tilde{\alpha}_4 \tilde{\alpha}_6} \right). \end{aligned} \tag{\infty.107}$$

This tensor is a generalization of the familiar inverting matrix $X_I^{\tilde{\alpha}_1 \tilde{\alpha}_2} \equiv \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_2}$ for geometric sums over matrices that satisfies

$$\delta_{\tilde{\alpha}_3}^{\tilde{\alpha}_1} = \sum_{\tilde{\alpha}_2 \in \mathcal{A}} X_I^{\tilde{\alpha}_1 \tilde{\alpha}_2} \frac{1}{\eta} \left[\delta_{\tilde{\alpha}_2 \tilde{\alpha}_3} - (\delta_{\tilde{\alpha}_2 \tilde{\alpha}_3} - \eta \tilde{H}_{\tilde{\alpha}_2 \tilde{\alpha}_3}) \right] = \sum_{\tilde{\alpha}_2 \in \mathcal{A}} X_I^{\tilde{\alpha}_1 \tilde{\alpha}_2} \tilde{H}_{\tilde{\alpha}_2 \tilde{\alpha}_3} \tag{\infty.108}$$

and appeared in the last expression of the dynamical helper function $a_{j;\tilde{\alpha}}(t)$ (∞.92). While we are on fire like an activated neuron, let's also define a final dynamical helper function for sums that are cubic in the free residual training error:

$$\begin{aligned}
 c_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}(t) &\equiv \eta \sum_{s=0}^t \epsilon_{j_1; \tilde{\alpha}_1}^{\text{F}}(t) \epsilon_{j_2; \tilde{\alpha}_2}^{\text{F}}(t) \epsilon_{j_3; \tilde{\alpha}_3}^{\text{F}}(t) \\
 &= \sum_{\tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6} X_{\text{III}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \left[(z_{k_1; \tilde{\alpha}_4} - y_{k_1; \tilde{\alpha}_4}) (z_{k_2; \tilde{\alpha}_5} - y_{k_2; \tilde{\alpha}_5}) (z_{k_3; \tilde{\alpha}_6} - y_{k_3; \tilde{\alpha}_6}) \right. \\
 &\quad \left. - \epsilon_{j_1; \tilde{\alpha}_4}^{\text{F}}(t) \epsilon_{j_2; \tilde{\alpha}_5}^{\text{F}}(t) \epsilon_{j_3; \tilde{\alpha}_6}^{\text{F}}(t) \right] + O\left(\frac{1}{n}\right).
 \end{aligned}
 \tag{\infty.109}$$

In this case, the inverting tensor $X_{\text{III}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}$ is implicitly defined by

$$\begin{aligned}
 \delta_{\tilde{\alpha}_7}^{\tilde{\alpha}_1} \delta_{\tilde{\alpha}_8}^{\tilde{\alpha}_2} \delta_{\tilde{\alpha}_9}^{\tilde{\alpha}_3} &= \sum_{\tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6} X_{\text{III}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \left[\left(\tilde{H}_{\tilde{\alpha}_4 \tilde{\alpha}_7} \delta_{\tilde{\alpha}_5 \tilde{\alpha}_8} \delta_{\tilde{\alpha}_6 \tilde{\alpha}_9} + \delta_{\tilde{\alpha}_4 \tilde{\alpha}_7} \tilde{H}_{\tilde{\alpha}_5 \tilde{\alpha}_8} \delta_{\tilde{\alpha}_6 \tilde{\alpha}_9} + \delta_{\tilde{\alpha}_4 \tilde{\alpha}_7} \delta_{\tilde{\alpha}_5 \tilde{\alpha}_8} \tilde{H}_{\tilde{\alpha}_6 \tilde{\alpha}_9} \right) \right. \\
 &\quad \left. - \eta \left(\tilde{H}_{\tilde{\alpha}_4 \tilde{\alpha}_7} \tilde{H}_{\tilde{\alpha}_5 \tilde{\alpha}_8} \delta_{\tilde{\alpha}_6 \tilde{\alpha}_9} + \tilde{H}_{\tilde{\alpha}_4 \tilde{\alpha}_7} \delta_{\tilde{\alpha}_5 \tilde{\alpha}_8} \tilde{H}_{\tilde{\alpha}_6 \tilde{\alpha}_9} + \delta_{\tilde{\alpha}_4 \tilde{\alpha}_7} \tilde{H}_{\tilde{\alpha}_5 \tilde{\alpha}_8} \tilde{H}_{\tilde{\alpha}_6 \tilde{\alpha}_9} \right) \right. \\
 &\quad \left. + \eta^2 \left(\tilde{H}_{\tilde{\alpha}_4 \tilde{\alpha}_7} \tilde{H}_{\tilde{\alpha}_5 \tilde{\alpha}_8} \tilde{H}_{\tilde{\alpha}_6 \tilde{\alpha}_9} \right) \right].
 \end{aligned}
 \tag{\infty.110}$$

Essentially, these three dynamical helper functions encode the step dependence of various geometric sums of the free residual training error $\epsilon_{j; \tilde{\alpha}}^{\text{F}}(t)$. Note that we have

$$a_{j_1; \tilde{\alpha}_1}(t=0) = 0, \tag{\infty.111}$$

$$b_{j_1 j_2; \tilde{\alpha}_1 \tilde{\alpha}_2}(t=0) = 0, \tag{\infty.112}$$

$$c_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}(t=0) = 0, \tag{\infty.113}$$

and so they all vanish at initialization, while at the end of training we have

$$a_{j_1; \tilde{\alpha}_1}(\infty) = \sum_{\tilde{\alpha}_2} X_{\text{I}}^{\tilde{\alpha}_1 \tilde{\alpha}_2} (z_{j_1; \tilde{\alpha}_2} - y_{j_1; \tilde{\alpha}_2}) + O\left(\frac{1}{n}\right), \tag{\infty.114}$$

$$b_{j_1 j_2; \tilde{\alpha}_1 \tilde{\alpha}_2}(\infty) = \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4} X_{\text{II}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (z_{j_1; \tilde{\alpha}_3} - y_{j_1; \tilde{\alpha}_3}) (z_{j_2; \tilde{\alpha}_4} - y_{j_2; \tilde{\alpha}_4}) + O\left(\frac{1}{n}\right), \tag{\infty.115}$$

$$\begin{aligned}
 c_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}(\infty) &= \sum_{\tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6} X_{\text{III}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) \\
 &\quad + O\left(\frac{1}{n}\right),
 \end{aligned}
 \tag{\infty.116}$$

since the free residual training error $\epsilon_{j; \tilde{\alpha}}^{\text{F}}(t)$ vanishes exponentially quickly, cf. (∞.90).

With the help of the first two of these dynamical helper functions, we can semi-compactly write the solution to the difference equation for the interaction NTK (∞.105) as

$$\begin{aligned}
 & H_{i_1 i_2; \delta_1 \delta_2}^I(t) \\
 &= - \sum_{j, \tilde{\alpha}} \left(\widehat{dH}_{i_1 i_2 j; \delta_1 \delta_2 \tilde{\alpha}} + \widehat{dH}_{i_2 i_1 j; \delta_2 \delta_1 \tilde{\alpha}} \right) a_{j; \tilde{\alpha}}(t) \\
 &+ \sum_{j, k, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3} \left(\widehat{dd_I H}_{i_1 i_2 j k; \delta_1 \delta_2 \tilde{\alpha}_1 \tilde{\alpha}_2} + \widehat{dd_{II} H}_{i_1 i_2 j k; \delta_1 \delta_2 \tilde{\alpha}_1 \tilde{\alpha}_2} + \widehat{dd_{II} H}_{i_1 j i_2 k; \delta_1 \tilde{\alpha}_1 \delta_2 \tilde{\alpha}_2} \right. \\
 &\quad \left. + \widehat{dd_I H}_{i_2 i_1 j k; \delta_2 \delta_1 \tilde{\alpha}_1 \tilde{\alpha}_2} + \widehat{dd_{II} H}_{i_2 i_1 j k; \delta_2 \delta_1 \tilde{\alpha}_1 \tilde{\alpha}_2} + \widehat{dd_{II} H}_{i_2 j i_1 k; \delta_2 \tilde{\alpha}_1 \delta_1 \tilde{\alpha}_2} \right) \\
 &\quad \times \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_3} \left[(z_{k; \tilde{\alpha}_3} - y_{k; \tilde{\alpha}_3}) a_{j; \tilde{\alpha}_1}(t) - b_{j k; \tilde{\alpha}_1 \tilde{\alpha}_3}(t) \right] \\
 &+ \frac{\eta}{2} \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2} \left(\widehat{dd_I H}_{i_1 i_2 j_1 j_2; \delta_1 \delta_2 \tilde{\alpha}_1 \tilde{\alpha}_2} + \widehat{dd_I H}_{i_2 i_1 j_1 j_2; \delta_2 \delta_1 \tilde{\alpha}_1 \tilde{\alpha}_2} \right. \\
 &\quad \left. + 2 \widehat{dd_{II} H}_{i_1 i_2 j_1 j_2; \delta_1 \delta_2 \tilde{\alpha}_1 \tilde{\alpha}_2} \right) b_{j_1 j_2; \tilde{\alpha}_1 \tilde{\alpha}_2}(t).
 \end{aligned} \tag{\infty.117}$$

As a quick sanity check, the vanishing initial condition for the interaction NTK, (∞.99), is satisfied due to the vanishing of the helper functions at initialization, (∞.111) and (∞.112). We can also plug in (∞.114) and (∞.115) to evaluate the change in NTK at the end of training. In particular, we see that the larger the change in the NTK, the more the feature functions evolve. Thus, more initial error entails more representation learning over the course of training.

Lastly, we need to determine the step dependence of the interacting part of the network output, $z_{i; \delta}^I(t)$. Inserting our free-interacting decomposition, (∞.94), into our dynamics for the network output, (∞.81), and using the fact that the free part satisfies the step-independent evolution equation (∞.86), we can reorganize terms to find a dynamical equation for the interacting part only:

$$z_{i; \delta}^I(t+1) = z_{i; \delta}^I(t) - \sum_{j, \tilde{\alpha}} \eta \widehat{H}_{ij; \delta \tilde{\alpha}} z_{j; \tilde{\alpha}}^I(t) + \eta \mathbb{F}_{i; \delta}(t). \tag{\infty.118}$$

Here, we've defined a *damping force*:

$$\begin{aligned}
 \mathbb{F}_{i; \delta}(t) \equiv & - \sum_{j, \tilde{\alpha}} H_{ij; \delta \tilde{\alpha}}^I(t) \epsilon_{j; \tilde{\alpha}}^F(t) + \frac{\eta}{2} \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2} dH_{ij_1 j_2; \delta \tilde{\alpha}_1 \tilde{\alpha}_2}(t) \epsilon_{j_1; \tilde{\alpha}_1}^F(t) \epsilon_{j_2; \tilde{\alpha}_2}^F(t) \\
 & - \frac{\eta^2}{6} \sum_{j_1, j_2, j_3, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3} \widehat{dd_I H}_{ij_1 j_2 j_3; \delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \epsilon_{j_1; \tilde{\alpha}_1}^F(t) \epsilon_{j_2; \tilde{\alpha}_2}^F(t) \epsilon_{j_3; \tilde{\alpha}_3}^F(t).
 \end{aligned} \tag{\infty.119}$$

Since we have solutions for the interaction NTK and the dNTK dynamics in terms of the free residual training error solution $\epsilon_{j; \tilde{\alpha}}^F(t)$, (∞.117) and (∞.102), this damping force is an explicitly known function of the step t .

Let us now try to implicitly express the solution to this difference equation (∞.118) as a sum over steps. First, for inputs in the training set $\tilde{\alpha} \in \mathcal{A}$, the dynamics of the interacting part of the output, (∞.118), reduce to a first-order inhomogeneous linear difference equation, with the damping force ruining the homogeneity:

$$z_{i; \tilde{\alpha}}^I(t+1) = \sum_{j, \tilde{\alpha}_1} \left(I_{ij; \tilde{\alpha} \tilde{\alpha}_1} - \eta \widehat{H}_{ij; \tilde{\alpha} \tilde{\alpha}_1} \right) z_{j; \tilde{\alpha}_1}^I(t) + \eta \mathbb{F}_{i; \tilde{\alpha}}(t). \tag{\infty.120}$$

We can formally solve this equation as a convolution of the damping force with the free step-evolution operator (∞.89),

$$z_{i;\tilde{\alpha}}^I(t) = \eta \sum_{s=0}^{t-1} \sum_{j,\tilde{\alpha}_1} U_{ij;\tilde{\alpha}\tilde{\alpha}_1}(t-1-s) \mathbb{F}_{j;\tilde{\alpha}_1}(s), \quad (\infty.121)$$

which satisfies the initial condition $z_{i;\tilde{\alpha}}^I(t=0) = 0$. Plugging this result back into the dynamical equation for general inputs $\delta \in \mathcal{D}$, (∞.118), we can then find a solution for the interacting part of the network output for such general inputs:

$$z_{i;\delta}^I(t) = \eta \sum_{s=0}^{t-1} \left[\mathbb{F}_{i;\delta}(s) - \sum_{j,\tilde{\alpha}} \hat{H}_{ij;\delta\tilde{\alpha}} z_{j;\tilde{\alpha}}^I(s) \right]. \quad (\infty.122)$$

In order to further simplify these expressions, we need to work out one convoluted sum:

$$\begin{aligned} \eta \sum_{s=0}^{t-1} z_{j;\tilde{\alpha}}^I(s) &= \eta^2 \sum_{s=0}^{t-1} \sum_{\tilde{s}=0}^{s-1} \sum_{k,\tilde{\alpha}_1} U_{jk;\tilde{\alpha}\tilde{\alpha}_1}(s-1-\tilde{s}) \mathbb{F}_{k;\tilde{\alpha}_1}(\tilde{s}) \\ &= \eta \sum_{u=0}^{t-2} \sum_{k,\tilde{\alpha}_1} \left[\eta \sum_{\tilde{u}=0}^{t-u-2} U_{jk;\tilde{\alpha}\tilde{\alpha}_1}(\tilde{u}) \right] \mathbb{F}_{k;\tilde{\alpha}_1}(u) \\ &= \eta \sum_{u=0}^{t-2} \sum_{k,m,\tilde{\alpha}_1,\tilde{\alpha}_2} \left(\hat{H}^{-1} \right)_{jm}^{\tilde{\alpha}\tilde{\alpha}_2} [I - U(t-u-1)]_{mk;\tilde{\alpha}_2\tilde{\alpha}_1} \mathbb{F}_{k;\tilde{\alpha}_1}(u) \\ &= \sum_{m,\tilde{\alpha}_2} \left(\hat{H}^{-1} \right)_{jm}^{\tilde{\alpha}\tilde{\alpha}_2} \left\{ \left[\eta \sum_{u=0}^{t-2} \mathbb{F}_{m;\tilde{\alpha}_2}(u) \right] + \left[\eta \mathbb{F}_{m;\tilde{\alpha}_2}(t-1) - z_{m;\tilde{\alpha}_2}^I(t) \right] \right\} \\ &= \sum_{m,\tilde{\alpha}_2} \left(\hat{H}^{-1} \right)_{jm}^{\tilde{\alpha}\tilde{\alpha}_2} \left\{ \left[\eta \sum_{u=0}^{t-1} \mathbb{F}_{m;\tilde{\alpha}_2}(u) \right] - z_{m;\tilde{\alpha}_2}^I(t) \right\}. \end{aligned} \quad (\infty.123)$$

Step by step, on the first line we used the expression for the formal solution (∞.121); on the second line we first reversed the order of the sums as $\sum_{s=0}^{t-1} \sum_{\tilde{s}=0}^{s-1} = \sum_{\tilde{s}=0}^{t-2} \sum_{s=\tilde{s}+1}^{t-1}$, and then we rewrote these sums in terms of new variables, $u \equiv \tilde{s}$ and $\tilde{u} \equiv s-1-\tilde{s}$; on the third line, we performed the geometric sum exactly as in (∞.92); on the fourth line, we used our formal solution (∞.121) at step t ; and on the final line, we combined terms to extend the limits of the sum over the damping force. Plugging this evaluation back into our formal solution for $z_{i;\delta}^I(t)$, (∞.122), we get a slightly less formal solution:

$$z_{i;\delta}^I(t) = \eta \sum_{s=0}^{t-1} \left[\mathbb{F}_{i;\delta}(s) - \sum_{\tilde{\alpha}_1,\tilde{\alpha}_2} H_{\delta\tilde{\alpha}_1} \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} \mathbb{F}_{i;\tilde{\alpha}_2}(s) \right] + \sum_{\tilde{\alpha}_1,\tilde{\alpha}_2} H_{\delta\tilde{\alpha}_1} \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} z_{i;\tilde{\alpha}_2}^I(t) + O\left(\frac{1}{n^2}\right). \quad (\infty.124)$$

In this result, we've replaced the stochastic NTK by its mean, $\hat{H}_{ij;\delta_1\delta_2} \rightarrow \delta_{ij} H_{\delta_1,\delta_2}$, as every term is otherwise already proportional to the dNTK or ddNTKs. As a quick sanity

check, note that for inputs in the training set, $\delta \rightarrow \tilde{\alpha} \in \mathcal{A}$, this general expression reduces to an identity on $z_{i;\tilde{\alpha}}^{\text{I}}(t)$.

Ultimately, what we care about is the interacting solution at the end of training, $t \rightarrow \infty$. As before, let's assume that the product of η with the stochastic NTK is sufficiently small such that the free step-evolution operator,

$$\lim_{t \rightarrow \infty} U(t) \propto \exp(-\eta \hat{H}t), \quad (\infty.125)$$

exponentially decays to zero.¹⁹ Then, for training inputs, the interacting part of the network outputs, $z_{i;\tilde{\alpha}}^{\text{I}}(t)$, will exponentially converge to zero:

$$\lim_{t \rightarrow \infty} z_{i;\tilde{\alpha}}^{\text{I}}(t) = 0. \quad (\infty.126)$$

To see why this holds, note that the dampening force ($\infty.119$) decays exponentially as

$$\lim_{s \rightarrow \infty} \mathbb{F}(s) \propto \exp(-\eta \hat{H}s), \quad (\infty.127)$$

since its leading behavior is linearly proportional to the free residual training error $\epsilon_{j;\tilde{\alpha}}^{\text{F}}(t)$. Combined with ($\infty.125$), this means that the interacting solution ($\infty.121$) converges as

$$\begin{aligned} \lim_{t \rightarrow \infty} z_{i;\tilde{\alpha}}^{\text{I}}(t) &= \eta \lim_{t \rightarrow \infty} \sum_{s=0}^{t-1} \sum_{j,\tilde{\alpha}_1} U_{ij;\tilde{\alpha}\tilde{\alpha}_1}(t-1-s) \mathbb{F}_{j;\tilde{\alpha}_1}(s) \\ &\propto \lim_{t \rightarrow \infty} \left\{ \eta t \exp[-(t-1)\eta \hat{H}] \right\} = 0, \end{aligned} \quad (\infty.128)$$

which is slightly slower than the convergence of the free solution $z_{i;\tilde{\alpha}}^{\text{F}}(t) \propto \exp(-\hat{H}t)$. Thus, overall the training algorithm converges:

$$\lim_{t \rightarrow \infty} z_{i;\tilde{\alpha}}(t) - y_{i;\tilde{\alpha}} = \lim_{t \rightarrow \infty} \left[z_{i;\tilde{\alpha}}^{\text{F}}(t) + z_{i;\tilde{\alpha}}^{\text{I}}(t) \right] - y_{i;\tilde{\alpha}} = 0, \quad (\infty.129)$$

where here we also used the free solution ($\infty.90$).²⁰ Incidentally, and of possible broader interest, this altogether shows that gradient descent converges exponentially quickly to a zero-error minimum for realistic deep neural networks of finite width and nonzero depth, up to errors that are at most quadratic in our effective theory cutoff L/n .

For general inputs, the interacting part of the output, $z_{i;\delta}^{\text{I}}(t)$, is given by the expression ($\infty.124$). With the convergence on the training set in mind, ($\infty.126$), the expression in the end-of-training limit reduces to

$$\lim_{t \rightarrow \infty} z_{i;\delta}^{\text{I}}(t) = \left[\eta \sum_{s=0}^{\infty} \mathbb{F}_{i;\delta}(s) \right] - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\delta\tilde{\alpha}_1} \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} \left[\eta \sum_{s=0}^{\infty} \mathbb{F}_{i;\tilde{\alpha}_2}(s) \right]. \quad (\infty.130)$$

¹⁹Note that since the step-evolution operator $U(t)$ is constructed in ($\infty.89$) from the *step-independent* NTK, $\hat{H}_{ij;\delta\tilde{\alpha}}$, the condition for this convergence is the same as the free analysis discussed in footnote 16.

²⁰Remember that in this section we have stopped explicitly denoting that there are $O(1/n^2)$ corrections. Recalling our fully-trained condition ($\infty.66$), this result ($\infty.129$) should be understood to be true up to such corrections, cf. ($\infty.78$) for our two-step solution where the situation was analogous.

Thus, all that remains is for us to perform an infinite sum over the damping force:

$$\begin{aligned} \eta \sum_{s=0}^{\infty} \mathbb{F}_{i;\delta}(s) \equiv & - \sum_{j, \tilde{\alpha}_1} \left[\eta \sum_{s=0}^{\infty} H_{ij;\delta \tilde{\alpha}_1}^{\text{I}}(s) \epsilon_{j;\tilde{\alpha}_1}^{\text{F}}(s) \right] \\ & + \frac{\eta}{2} \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2} \left[\eta \sum_{s=0}^{\infty} dH_{ij_1 j_2; \delta \tilde{\alpha}_1 \tilde{\alpha}_2}(s) \epsilon_{j_1; \tilde{\alpha}_1}^{\text{F}}(s) \epsilon_{j_2; \tilde{\alpha}_2}^{\text{F}}(s) \right] \\ & - \frac{\eta^2}{6} \sum_{j_1, j_2, j_3, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3} \left[\eta \sum_{s=0}^{\infty} \widehat{\text{dd}}_{\text{I}} H_{ij_1 j_2 j_3; \delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \epsilon_{j_1; \tilde{\alpha}_1}^{\text{F}}(s) \epsilon_{j_2; \tilde{\alpha}_2}^{\text{F}}(s) \epsilon_{j_3; \tilde{\alpha}_3}^{\text{F}}(s) \right]. \end{aligned} \quad (\infty.131)$$

The third sum is exactly the end-of-training limit of the third dynamical helper function $c_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}(t)$ that we already evaluated in (∞.116), giving

$$\begin{aligned} & \eta \sum_{s=0}^{\infty} \widehat{\text{dd}}_{\text{I}} H_{ij_1 j_2 j_3; \delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \epsilon_{j_1; \tilde{\alpha}_1}^{\text{F}}(s) \epsilon_{j_2; \tilde{\alpha}_2}^{\text{F}}(s) \epsilon_{j_3; \tilde{\alpha}_3}^{\text{F}}(s) \\ & = \sum_{\tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6} \widehat{\text{dd}}_{\text{I}} H_{ij_1 j_2 j_3; \delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} X_{\text{III}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}). \end{aligned} \quad (\infty.132)$$

Thus, we are left with two more sums to evaluate. Let's proceed slowly: everything is simple, but there are a lot of terms to get right.

To start, we can evaluate the second sum in (∞.131) as

$$\begin{aligned} & \eta \sum_{s=0}^{\infty} \widehat{d}H_{ij_1 j_2; \delta \tilde{\alpha}_1 \tilde{\alpha}_2}(s) \epsilon_{j_1; \tilde{\alpha}_1}^{\text{F}}(s) \epsilon_{j_2; \tilde{\alpha}_2}^{\text{F}}(s) \\ & = \widehat{d}H_{ij_1 j_2; \delta \tilde{\alpha}_1 \tilde{\alpha}_2} \sum_{s=0}^{\infty} \left[\eta \epsilon_{j_1; \tilde{\alpha}_1}^{\text{F}}(s) \epsilon_{j_2; \tilde{\alpha}_2}^{\text{F}}(s) \right] \\ & \quad - \sum_{k, \tilde{\alpha}_3, \tilde{\alpha}_4} \left(\widehat{\text{dd}}_{\text{I}} H_{ij_1 j_2 k; \delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} + 2 \widehat{\text{dd}}_{\text{II}} H_{ij_1 j_2 k; \delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right) \\ & \quad \times \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_4} \sum_{s=0}^{\infty} \left\{ \eta \epsilon_{j_1; \tilde{\alpha}_1}^{\text{F}}(s) \epsilon_{j_2; \tilde{\alpha}_2}^{\text{F}}(s) \left[z_{k; \tilde{\alpha}_4} - y_{k; \tilde{\alpha}_4} - \epsilon_{k; \tilde{\alpha}_4}^{\text{F}}(s) \right] \right\}, \end{aligned} \quad (\infty.133)$$

where we substituted in our dynamical dNTK solution, (∞.102), and used the fact that the overall expression is symmetric under $(\tilde{\alpha}_1, j_1) \leftrightarrow (\tilde{\alpha}_2, j_2)$ to combine the two $\widehat{\text{dd}}_{\text{II}} H$ terms. Then, using our already evaluated sums, (∞.115) and (∞.116), we get

$$\begin{aligned} & \eta \sum_{s=0}^{\infty} \widehat{d}H_{ij_1 j_2; \delta \tilde{\alpha}_1 \tilde{\alpha}_2}(s) \epsilon_{j_1; \tilde{\alpha}_1}^{\text{F}}(s) \epsilon_{j_2; \tilde{\alpha}_2}^{\text{F}}(s) \\ & = \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4} \widehat{d}H_{ij_1 j_2; \delta \tilde{\alpha}_1 \tilde{\alpha}_2} X_{\text{II}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (z_{j_1; \tilde{\alpha}_3} - y_{j_1; \tilde{\alpha}_3}) (z_{j_2; \tilde{\alpha}_4} - y_{j_2; \tilde{\alpha}_4}) \\ & \quad - \sum_{k, \tilde{\alpha}_3, \dots, \tilde{\alpha}_6} \left(\widehat{\text{dd}}_{\text{I}} H_{ij_1 j_2 k; \delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} + 2 \widehat{\text{dd}}_{\text{II}} H_{ij_1 j_2 k; \delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right) \\ & \quad \times Y_1^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{k; \tilde{\alpha}_6} - y_{k; \tilde{\alpha}_6}), \end{aligned} \quad (\infty.134)$$

where we introduced a shorthand

$$Y_1^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \equiv X_{\text{II}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_4 \tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_6} - \sum_{\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_7} X_{\text{III}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_7 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}, \quad (\infty.135)$$

to ease the collection of terms later.

To finish, let us write the first sum in ($\infty.131$) as

$$\begin{aligned} & \eta \sum_{s=0}^{\infty} H_{ij;\delta\tilde{\alpha}_1}^{\text{I}}(s) \epsilon_{j;\tilde{\alpha}_1}^{\text{F}}(s) \\ &= - \sum_{k,\tilde{\alpha}_2} \left(\widehat{\text{d}H}_{ijk;\delta\tilde{\alpha}_1\tilde{\alpha}_2} + \widehat{\text{d}H}_{jik;\tilde{\alpha}_1\delta\tilde{\alpha}_2} \right) \sum_{s=0}^{\infty} \left[\eta \epsilon_{j;\tilde{\alpha}_1}^{\text{F}}(s) a_{k;\tilde{\alpha}_2}(s) \right] \\ &+ \sum_{k_1,k_2,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4} \left(\widehat{\text{dd}_{\text{I}}H}_{ijk_1k_2;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} + \widehat{\text{dd}_{\text{II}}H}_{ijk_1k_2;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} + \widehat{\text{dd}_{\text{II}}H}_{ik_1jk_2;\delta\tilde{\alpha}_2\tilde{\alpha}_1\tilde{\alpha}_3} \right. \\ &\quad \left. + \widehat{\text{dd}_{\text{I}}H}_{jik_1k_2;\tilde{\alpha}_1\delta\tilde{\alpha}_2\tilde{\alpha}_3} + \widehat{\text{dd}_{\text{II}}H}_{jik_1k_2;\tilde{\alpha}_1\delta\tilde{\alpha}_2\tilde{\alpha}_3} + \widehat{\text{dd}_{\text{II}}H}_{jk_1ik_2;\tilde{\alpha}_1\tilde{\alpha}_2\delta\tilde{\alpha}_3} \right) \\ &\quad \times \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} \sum_{s=0}^{\infty} \left\{ \eta \epsilon_{j;\tilde{\alpha}_1}^{\text{F}}(s) \left[(z_{k_2;\tilde{\alpha}_4} - y_{k_2;\tilde{\alpha}_4}) a_{k_1;\tilde{\alpha}_2}(s) - b_{k_1k_2;\tilde{\alpha}_2\tilde{\alpha}_4}(s) \right] \right\} \\ &+ \frac{\eta}{2} \sum_{k_1,k_2,\tilde{\alpha}_2,\tilde{\alpha}_3} \left(\widehat{\text{dd}_{\text{I}}H}_{ijk_1k_2;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} + \widehat{\text{dd}_{\text{I}}H}_{jik_1k_2;\tilde{\alpha}_1\delta\tilde{\alpha}_2\tilde{\alpha}_3} + 2\widehat{\text{dd}_{\text{II}}H}_{ijk_1k_2;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} \right) \\ &\quad \times \sum_{s=0}^{\infty} \left[\eta \epsilon_{j;\tilde{\alpha}_1}^{\text{F}}(s) b_{k_1k_2;\tilde{\alpha}_2\tilde{\alpha}_3}(s) \right], \end{aligned} \quad (\infty.136)$$

where we substituted in our solution for the interaction NTK, ($\infty.117$). Then, substituting for the helper functions with ($\infty.92$) and ($\infty.106$), performing the *additional* sums over these terms, and then using the end-of-training limits ($\infty.114$)–($\infty.116$), we get

$$\begin{aligned} & \eta \sum_{s=0}^{\infty} H_{ij;\delta\tilde{\alpha}_1}^{\text{I}}(s) \epsilon_{j;\tilde{\alpha}_1}^{\text{F}}(s) \\ &= - \sum_{k,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4} \left(\widehat{\text{d}H}_{ijk;\delta\tilde{\alpha}_1\tilde{\alpha}_2} + \widehat{\text{d}H}_{jik;\tilde{\alpha}_1\delta\tilde{\alpha}_2} \right) Y_2^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4} (z_{j;\tilde{\alpha}_3} - y_{j;\tilde{\alpha}_3}) (z_{k;\tilde{\alpha}_4} - y_{k;\tilde{\alpha}_4}) \\ &+ \sum_{k_1,k_2,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6} \left(\widehat{\text{dd}_{\text{I}}H}_{ijk_1k_2;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} + \widehat{\text{dd}_{\text{II}}H}_{ijk_1k_2;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} + \widehat{\text{dd}_{\text{II}}H}_{ik_1jk_2;\delta\tilde{\alpha}_2\tilde{\alpha}_1\tilde{\alpha}_3} \right. \\ &\quad \left. + \widehat{\text{dd}_{\text{I}}H}_{jik_1k_2;\tilde{\alpha}_1\delta\tilde{\alpha}_2\tilde{\alpha}_3} + \widehat{\text{dd}_{\text{II}}H}_{jik_1k_2;\tilde{\alpha}_1\delta\tilde{\alpha}_2\tilde{\alpha}_3} + \widehat{\text{dd}_{\text{II}}H}_{jk_1ik_2;\tilde{\alpha}_1\tilde{\alpha}_2\delta\tilde{\alpha}_3} \right) \\ &\quad \times Y_3^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j;\tilde{\alpha}_4} - y_{j;\tilde{\alpha}_4}) (z_{k_1;\tilde{\alpha}_5} - y_{k_1;\tilde{\alpha}_5}) (z_{k_2;\tilde{\alpha}_6} - y_{k_2;\tilde{\alpha}_6}) \\ &+ \frac{\eta}{2} \sum_{k_1,k_2,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6} \left(\widehat{\text{dd}_{\text{I}}H}_{ijk_1k_2;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} + \widehat{\text{dd}_{\text{I}}H}_{jik_1k_2;\tilde{\alpha}_1\delta\tilde{\alpha}_2\tilde{\alpha}_3} + 2\widehat{\text{dd}_{\text{II}}H}_{ijk_1k_2;\delta\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3} \right) \\ &\quad \times Y_4^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4\tilde{\alpha}_5\tilde{\alpha}_6} (z_{j;\tilde{\alpha}_4} - y_{j;\tilde{\alpha}_4}) (z_{k_1;\tilde{\alpha}_5} - y_{k_1;\tilde{\alpha}_5}) (z_{k_2;\tilde{\alpha}_6} - y_{k_2;\tilde{\alpha}_6}), \end{aligned} \quad (\infty.137)$$

where we introduced additional shorthands

$$Y_2^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \equiv \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_4} - \sum_{\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_5} X_{\text{II}}^{\tilde{\alpha}_1 \tilde{\alpha}_5 \tilde{\alpha}_3 \tilde{\alpha}_4}, \quad (\infty.138)$$

$$Y_3^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \equiv \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_4} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_6} - \sum_{\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_7} X_{\text{II}}^{\tilde{\alpha}_1 \tilde{\alpha}_7 \tilde{\alpha}_4 \tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_6} \\ - \sum_{\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_4} \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_7} X_{\text{II}}^{\tilde{\alpha}_2 \tilde{\alpha}_7 \tilde{\alpha}_5 \tilde{\alpha}_6} + \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8, \tilde{\alpha}_9} \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_9} X_{\text{II}}^{\tilde{\alpha}_2 \tilde{\alpha}_9 \tilde{\alpha}_7 \tilde{\alpha}_8} X_{\text{III}}^{\tilde{\alpha}_1 \tilde{\alpha}_7 \tilde{\alpha}_8 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}, \quad (\infty.139)$$

$$Y_4^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \equiv \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_4} X_{\text{II}}^{\tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_5 \tilde{\alpha}_6} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} X_{\text{II}}^{\tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_7 \tilde{\alpha}_8} X_{\text{III}}^{\tilde{\alpha}_1 \tilde{\alpha}_7 \tilde{\alpha}_8 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}. \quad (\infty.140)$$

Now, it's time to frantically flip through pages and collect everything we computed. Plugging the sums (∞.137), (∞.134), and (∞.132) back into our expression for the sum over the damping force, (∞.131), plugging that back into our expression for the interaction part of the network output, (∞.130), and then combining with the free part of the network output, (∞.93), we obtain our fully-trained solution for finite-width networks trained by gradient descent:

$$z_{i;\delta}(t = \infty) \quad (\infty.141) \\ \equiv z_{i;\delta}^{\text{F}}(t = \infty) + z_{i;\delta}^{\text{I}}(t = \infty) \\ = z_{i;\delta} - \sum_{j,k,\tilde{\alpha}_1,\tilde{\alpha}_2} \hat{H}_{ij;\delta\tilde{\alpha}_1} \left(\hat{H}^{-1} \right)_{jk}^{\tilde{\alpha}_1 \tilde{\alpha}_2} (z_{k;\tilde{\alpha}_2} - y_{k;\tilde{\alpha}_2}) \\ + \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4} \left[\widehat{\text{dH}}_{j_1 i j_2; \tilde{\alpha}_1 \delta \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\delta \tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5 \tilde{\alpha}_6} \widehat{\text{dH}}_{j_1 i j_2; \tilde{\alpha}_1 \tilde{\alpha}_6 \tilde{\alpha}_2} \right] \\ \times Z_{\text{A}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (z_{j_1; \tilde{\alpha}_3} - y_{j_1; \tilde{\alpha}_3}) (z_{j_2; \tilde{\alpha}_4} - y_{j_2; \tilde{\alpha}_4}) \\ + \sum_{j_1,j_2,\tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4} \left[\widehat{\text{dH}}_{i j_1 j_2; \delta \tilde{\alpha}_1 \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\delta \tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5 \tilde{\alpha}_6} \widehat{\text{dH}}_{i j_1 j_2; \tilde{\alpha}_6 \tilde{\alpha}_1 \tilde{\alpha}_2} \right] \\ \times Z_{\text{B}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (z_{j_1; \tilde{\alpha}_3} - y_{j_1; \tilde{\alpha}_3}) (z_{j_2; \tilde{\alpha}_4} - y_{j_2; \tilde{\alpha}_4}) \\ + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[\widehat{\text{dd}_1 \text{H}}_{j_1 i j_2 j_3; \tilde{\alpha}_1 \delta \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\delta \tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7 \tilde{\alpha}_8} \widehat{\text{dd}_1 \text{H}}_{j_1 i j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_8 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] \\ \times Z_{\text{IA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) \\ + \sum_{\substack{j_1,j_2,j_3, \\ \tilde{\alpha}_1,\tilde{\alpha}_2,\tilde{\alpha}_3,\tilde{\alpha}_4,\tilde{\alpha}_5,\tilde{\alpha}_6}} \left[\widehat{\text{dd}_1 \text{H}}_{i j_1 j_2 j_3; \delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\delta \tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7 \tilde{\alpha}_8} \widehat{\text{dd}_1 \text{H}}_{i j_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] \\ \times Z_{\text{IB}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6})$$

$$\begin{aligned}
& + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{\text{dd}_{\text{II}} H}_{j_1 j_2 i j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \delta \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\delta \tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7 \tilde{\alpha}_8} \widehat{\text{dd}_{\text{II}} H}_{j_1 j_2 i j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_8 \tilde{\alpha}_3} \right] \\
& \quad \times Z_{\text{IIA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) \\
& + \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{\text{dd}_{\text{II}} H}_{i j_1 j_2 j_3; \delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\delta \tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7 \tilde{\alpha}_8} \widehat{\text{dd}_{\text{II}} H}_{i j_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] \\
& \quad \times Z_{\text{IIB}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) \\
& + O\left(\frac{1}{n^2}\right).
\end{aligned}$$

Here we defined our final tensors with our final alphabet letter (and various subscripts),

$$Z_{\text{A}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \equiv Y_2^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4}, \quad (\infty.142)$$

$$Z_{\text{B}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \equiv Y_2^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} + \frac{\eta}{2} X_{\text{II}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4}, \quad (\infty.143)$$

$$Z_{\text{IA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \equiv -Y_3^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} - \frac{\eta}{2} Y_4^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}, \quad (\infty.144)$$

$$\begin{aligned}
Z_{\text{IB}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} & \equiv -Y_3^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} - \frac{\eta}{2} Y_4^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \\
& \quad - \frac{\eta}{2} Y_1^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} - \frac{\eta^2}{6} X_{\text{III}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}, \quad (\infty.145)
\end{aligned}$$

$$Z_{\text{IIA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \equiv -Y_3^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}, \quad (\infty.146)$$

$$\begin{aligned}
Z_{\text{IIB}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} & \equiv -Y_3^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} - Y_3^{\tilde{\alpha}_2 \tilde{\alpha}_1 \tilde{\alpha}_3 \tilde{\alpha}_5 \tilde{\alpha}_4 \tilde{\alpha}_6} - Y_3^{\tilde{\alpha}_1 \tilde{\alpha}_3 \tilde{\alpha}_2 \tilde{\alpha}_4 \tilde{\alpha}_6 \tilde{\alpha}_5} \\
& \quad - \eta Y_4^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} - \eta Y_1^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}, \quad (\infty.147)
\end{aligned}$$

making use of our previous shorthand tensors ($\infty.135$) and ($\infty.138$)–($\infty.140$).²¹ These **algorithm projectors**, ($\infty.142$)–($\infty.147$), serve to project the initial training error onto two different combinations of the dNTK, two different combinations of the first ddNTK, and two other different combinations of the second ddNTK, all according to the details of the gradient descent algorithm. As a final sanity check, note that as for inputs in the training set, $\delta \rightarrow \tilde{\alpha} \in \mathcal{A}$, the quantities in the square brackets in the finite-width solution ($\infty.141$) each vanish, and we recover our fully-trained condition ($\infty.66$).

Before we retire this subsection, let us elaborate on the algorithm dependence. First, note that our two-update solution ($\infty.79$) has the same form as the gradient-descent solution ($\infty.141$), but with different algorithm projectors:

$$Z_{\text{B}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \equiv \frac{1}{2} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_4}, \quad Z_{\text{IB}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \equiv -\frac{1}{6} \tilde{H}^{\tilde{\alpha}_1 \tilde{\alpha}_4} \tilde{H}^{\tilde{\alpha}_2 \tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_3 \tilde{\alpha}_6}, \quad (\infty.148)$$

and all the others vanishing. Clearly these algorithms have very different inductive biases! Second, we can study the ODE limit of the dynamics by taking $\eta \rightarrow 0$, cf. footnote 17:

²¹Note that for the last of these tensors, ($\infty.147$), in order to coax the various contributions in our solution ($\infty.141$) into the proper form, we used the symmetry of $\widehat{\text{dd}_{\text{II}} H}_{i j_1 j_2 j_3; \delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}$ and relabeled various dummy sample indices.

in this case, we see that the ODE dynamics have a solution given by

$$Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} = Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \equiv Y_2^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4}, \quad (\infty.149)$$

$$Z_{IA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} = Z_{IB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} = Z_{IIA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \equiv -Y_3^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \quad (\infty.150)$$

$$Z_{IIB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \equiv -Y_3^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} - Y_3^{\tilde{\alpha}_2 \tilde{\alpha}_1 \tilde{\alpha}_3 \tilde{\alpha}_5 \tilde{\alpha}_4 \tilde{\alpha}_6} - Y_3^{\tilde{\alpha}_1 \tilde{\alpha}_3 \tilde{\alpha}_2 \tilde{\alpha}_4 \tilde{\alpha}_6 \tilde{\alpha}_5}, \quad (\infty.151)$$

where $Y_2^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4}$ and $Y_3^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}$ are given by (∞.138) and (∞.139), respectively, the inverting tensor $X_{II}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4}$, (∞.107), now satisfies

$$\sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} X_{II}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \left(\tilde{H}_{\tilde{\alpha}_3 \tilde{\alpha}_5} \delta_{\tilde{\alpha}_4 \tilde{\alpha}_6} + \delta_{\tilde{\alpha}_3 \tilde{\alpha}_5} \tilde{H}_{\tilde{\alpha}_4 \tilde{\alpha}_6} \right) = \delta_{\tilde{\alpha}_5}^{\tilde{\alpha}_1} \delta_{\tilde{\alpha}_6}^{\tilde{\alpha}_2}, \quad (\infty.152)$$

and the other inverting tensor $X_{II}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}$, (∞.110), now satisfies

$$\begin{aligned} & \delta_{\tilde{\alpha}_7}^{\tilde{\alpha}_1} \delta_{\tilde{\alpha}_8}^{\tilde{\alpha}_2} \delta_{\tilde{\alpha}_9}^{\tilde{\alpha}_3} \\ &= \sum_{\tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6} X_{III}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} \left(\tilde{H}_{\tilde{\alpha}_4 \tilde{\alpha}_7} \delta_{\tilde{\alpha}_5 \tilde{\alpha}_8} \delta_{\tilde{\alpha}_6 \tilde{\alpha}_9} + \delta_{\tilde{\alpha}_4 \tilde{\alpha}_7} \tilde{H}_{\tilde{\alpha}_5 \tilde{\alpha}_8} \delta_{\tilde{\alpha}_6 \tilde{\alpha}_9} + \delta_{\tilde{\alpha}_4 \tilde{\alpha}_7} \delta_{\tilde{\alpha}_5 \tilde{\alpha}_8} \tilde{H}_{\tilde{\alpha}_6 \tilde{\alpha}_9} \right). \end{aligned} \quad (\infty.153)$$

This entirely captures the difference between gradient flow and gradient descent for these fully-trained networks. In general, we conjecture that for finite-width networks, at leading order the fully-trained solution takes the universal form of (∞.141), with all of the *algorithm dependence* encoded by the six *algorithm projectors*: $Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4}$, $Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4}$, $Z_{IA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}$, $Z_{IB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}$, $Z_{IIA}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}$, and $Z_{IIB}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6}$.²²

Finally, let's understand what is meant by the remaining error in our solution (∞.141). It is not really the error of an actual network that is instantiated and then fully trained through many many gradient-descent steps, but instead it is the error in our effective description of such a particular fully-trained network. Of course, our effective theory formalism can compute higher-order corrections, if they're of interest. However, the leading-order finite-width corrections should really be sufficient for most cases: for instance, for a network of depth $L = 10$ layers and of hidden-layer width $n = 100$ neurons each, our effective description will only be off by $\sim (L/n)^2 = 1\%$.

Furthermore, for any theoretical analysis, the main qualitative difference in the solution appears when going from infinite width to finite width, as we go from a free theory to an interacting theory and from linear dynamics to nonlinear dynamics. Thus, the effective theory that gave us the solution (∞.141) really is “as simple... as possible” while still providing an extremely accurate description of real deep learning models.

∞.2.3 Prediction at Finite Width

Having solved the training dynamics in two different ways, we can now rather generally study the predictions of our networks on novel inputs $x_{\hat{\beta}}$ from the test set $\hat{\beta} \in \mathcal{B}$.

²²More precisely, we conjecture that our conjecture, as stated, holds for the MSE loss. For the cross-entropy loss the solution will take a slightly different form but with a similar partition into an algorithm-independent part and an algorithm-dependent part described by similar algorithm projectors.

At finite width, the predictions of a fully-trained network are universally governed by the stochastic equation

$$\begin{aligned}
 z_{i;\dot{\beta}}(t=T) &= z_{i;\dot{\beta}} - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\dot{\beta}\tilde{\alpha}_1} \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (z_{i;\tilde{\alpha}_2} - y_{i;\tilde{\alpha}_2}) \\
 &+ \sum_{j, \tilde{\alpha}_1, \tilde{\alpha}_2} \left[\widehat{\Delta H}_{ij;\dot{\beta}\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} H_{\dot{\beta}\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} \widehat{\Delta H}_{ij;\tilde{\alpha}_4\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (z_{j;\tilde{\alpha}_2} - y_{j;\tilde{\alpha}_2}) \\
 &- \sum_{\substack{j,k \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_4}} \left[\widehat{\Delta H}_{ij;\dot{\beta}\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\dot{\beta}\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{\Delta H}_{ij;\tilde{\alpha}_6\tilde{\alpha}_1} \right] \tilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} \widehat{\Delta H}_{jk;\tilde{\alpha}_2\tilde{\alpha}_3} \tilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} (z_{k;\tilde{\alpha}_4} - y_{k;\tilde{\alpha}_4}) \\
 &+ \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{\mathrm{d}H}_{j_1 j_2; \tilde{\alpha}_1 \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\dot{\beta}\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{\mathrm{d}H}_{j_1 j_2; \tilde{\alpha}_1 \tilde{\alpha}_6 \tilde{\alpha}_2} \right] \\
 &\quad \times Z_{\tilde{\mathbf{A}}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (z_{j_1; \tilde{\alpha}_3} - y_{j_1; \tilde{\alpha}_3}) (z_{j_2; \tilde{\alpha}_4} - y_{j_2; \tilde{\alpha}_4}) \\
 &+ \sum_{j_1, j_2, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \left[\widehat{\mathrm{d}H}_{ij_1 j_2; \dot{\beta} \tilde{\alpha}_1 \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6} H_{\dot{\beta}\tilde{\alpha}_5} \tilde{H}^{\tilde{\alpha}_5\tilde{\alpha}_6} \widehat{\mathrm{d}H}_{ij_1 j_2; \tilde{\alpha}_6 \tilde{\alpha}_1 \tilde{\alpha}_2} \right] \\
 &\quad \times Z_{\tilde{\mathbf{B}}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (z_{j_1; \tilde{\alpha}_3} - y_{j_1; \tilde{\alpha}_3}) (z_{j_2; \tilde{\alpha}_4} - y_{j_2; \tilde{\alpha}_4}) \\
 &+ \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{\mathrm{dd}_I H}_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\dot{\beta}\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{\mathrm{dd}_I H}_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_8 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] \\
 &\quad \times Z_{\tilde{\mathbf{I}}\tilde{\mathbf{A}}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) \\
 &+ \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{\mathrm{dd}_I H}_{ij_1 j_2 j_3; \dot{\beta} \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\dot{\beta}\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{\mathrm{dd}_I H}_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] \\
 &\quad \times Z_{\tilde{\mathbf{I}}\tilde{\mathbf{B}}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) \\
 &+ \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{\mathrm{dd}_{II} H}_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\dot{\beta}\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{\mathrm{dd}_{II} H}_{j_1 j_2 j_3; \tilde{\alpha}_1 \tilde{\alpha}_8 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] \\
 &\quad \times Z_{\tilde{\mathbf{II}}\tilde{\mathbf{A}}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) \\
 &+ \sum_{\substack{j_1, j_2, j_3, \\ \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4, \tilde{\alpha}_5, \tilde{\alpha}_6}} \left[\widehat{\mathrm{dd}_{II} H}_{ij_1 j_2 j_3; \dot{\beta} \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} - \sum_{\tilde{\alpha}_7, \tilde{\alpha}_8} H_{\dot{\beta}\tilde{\alpha}_7} \tilde{H}^{\tilde{\alpha}_7\tilde{\alpha}_8} \widehat{\mathrm{dd}_{II} H}_{ij_1 j_2 j_3; \tilde{\alpha}_8 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3} \right] \\
 &\quad \times Z_{\tilde{\mathbf{II}}\tilde{\mathbf{B}}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} (z_{j_1; \tilde{\alpha}_4} - y_{j_1; \tilde{\alpha}_4}) (z_{j_2; \tilde{\alpha}_5} - y_{j_2; \tilde{\alpha}_5}) (z_{j_3; \tilde{\alpha}_6} - y_{j_3; \tilde{\alpha}_6}) \\
 &+ O\left(\frac{1}{n^2}\right).
 \end{aligned}
 \tag{\infty.154}$$

This formula could borderline fit on a T-shirt.²³ For this expression, we've expanded the complete inverse of the stochastic NTK sub-tensor as per (∞.71); in particular, the second line is more or less the infinite-width kernel prediction (10.39), and the terms on the third and fourth lines are finite-width corrections due to NTK fluctuations. Further, the algorithm projectors (∞.142)–(∞.147) contain all the finite-width algorithm dependence of the solution, and thus this general solution (∞.154) can describe all our explicit solutions, whether we train in two steps (∞.79), in many many steps (∞.141), or with any other choice of optimization algorithm that uses the MSE loss.

²³To better facilitate such brand awareness, first recall that for nearly-kernel methods we were able to compress the model predictions in terms of a *trained kernel* (11.142); similarly, we can compress the predictions of finite-width networks (∞.154) in terms of a **trained NTK**, $H_{ij;\delta\tilde{\alpha}}^\sharp$, as

$$z_{i;\hat{\beta}}(t=T) = z_{i;\hat{\beta}} - \sum_{j,k,\tilde{\alpha}_1,\tilde{\alpha}_2} H_{ij;\hat{\beta}\tilde{\alpha}_1}^\sharp \widetilde{H}_{jk}^{\tilde{\alpha}_1\tilde{\alpha}_2} (z_{k;\tilde{\alpha}_2} - y_{k;\tilde{\alpha}_2}) + O\left(\frac{1}{n^2}\right), \quad (\infty.155)$$

taking the form of a (*neural tangent*) *kernel prediction* (10.39). To see how this works, let us decompose the trained NTK into free and training-dependent terms:

$$H_{ij;\delta\tilde{\alpha}}^\sharp \equiv \widehat{H}_{ij;\delta\tilde{\alpha}} + \mathbb{H}_{ij;\delta\tilde{\alpha}}. \quad (\infty.156)$$

Please don't confuse this decomposition with our earlier decomposition (∞.95): that former one was convenient for solving the training dynamics, while this new one is useful for determining the trained NTK in (∞.155). Considering the inverse of the trained NTK restricted to the training set only,

$$\widetilde{H}_{jk}^{\tilde{\alpha}_1\tilde{\alpha}_2} \equiv \left(\widehat{H}^{-1}\right)_{jk}^{\tilde{\alpha}_1\tilde{\alpha}_2} - \sum_{\tilde{\alpha}_3,\tilde{\alpha}_4} \mathbb{H}_{jk;\tilde{\alpha}_3\tilde{\alpha}_4} \widetilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_3} \widetilde{H}^{\tilde{\alpha}_2\tilde{\alpha}_4} + O\left(\frac{1}{n^2}\right), \quad (\infty.157)$$

and plugging it along with the decomposition (∞.156) into the formula (∞.155), we get

$$\begin{aligned} z_{i;\hat{\beta}}(t=T) = & z_{i;\hat{\beta}} - \sum_{j,k,\tilde{\alpha}_1,\tilde{\alpha}_2} \widehat{H}_{ij;\hat{\beta}\tilde{\alpha}_1} \left(\widehat{H}^{-1}\right)_{jk}^{\tilde{\alpha}_1\tilde{\alpha}_2} (z_{k;\tilde{\alpha}_2} - y_{k;\tilde{\alpha}_2}) \\ & - \sum_{j,\tilde{\alpha}_1,\tilde{\alpha}_2} \left[\mathbb{H}_{ij;\hat{\beta}\tilde{\alpha}_1} - \sum_{\tilde{\alpha}_3,\tilde{\alpha}_4} H_{\hat{\beta}\tilde{\alpha}_3} \widetilde{H}^{\tilde{\alpha}_3\tilde{\alpha}_4} \mathbb{H}_{ij;\tilde{\alpha}_4\tilde{\alpha}_1} \right] \widetilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (z_{j;\tilde{\alpha}_2} - y_{j;\tilde{\alpha}_2}) + O\left(\frac{1}{n^2}\right). \end{aligned} \quad (\infty.158)$$

The terms on the first line of the right-hand side of this expression give the free contribution to the solution, (∞.93), while the terms on the second line give the interacting contribution, (∞.130), encapsulating the effect of nontrivial representation learning at finite width. The specific form of $H_{ij;\delta\tilde{\alpha}}^\sharp$ can be found by matching the terms on the right-hand sides of (∞.158) and (∞.154). You can also express the training-dependent part, $\mathbb{H}_{ij;\delta\tilde{\alpha}}$, implicitly in terms of the damping force (∞.119) as

$$\eta \sum_{s=0}^{\infty} \mathbb{F}_{i;\delta}(s) = - \sum_{j,\tilde{\alpha}_1,\tilde{\alpha}_2} \mathbb{H}_{ij;\delta\tilde{\alpha}_1} \widetilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (z_{j;\tilde{\alpha}_2} - y_{j;\tilde{\alpha}_2}), \quad (\infty.159)$$

as is clear by comparing (∞.158) with (∞.130). This implicit expression also makes it clear that the form of $\mathbb{H}_{ij;\delta\tilde{\alpha}}$ isn't unique and can always be adjusted by the addition of a term orthogonal to $\sum_{\tilde{\alpha}_2} \widetilde{H}^{\tilde{\alpha}_1\tilde{\alpha}_2} (z_{j;\tilde{\alpha}_2} - y_{j;\tilde{\alpha}_2})$. In any event, with the trained NTK you can now fit the finite-width prediction formula (∞.155) on any newborn AI's onesie.

Furthermore, this solution (∞.154) describes the predictions of a *particular* network from our ensemble: the instantiation-to-instantiation difference is encoded in the particular initial network output, z , the NTK fluctuation, $\widehat{\Delta H}$, the dNTK, \widehat{dH} , the first ddNTK, $\widehat{dd_I H}$, and the second ddNTK, $\widehat{dd_{II} H}$. Since we know explicitly the statistics of these variables at initialization, we can also analyze the statistics of the fully-trained distribution in full. With an eye toward a discussion of the depth dependence of these statistics, we'll now revive the layer indices.

First and foremost, the mean prediction is given by

$$\begin{aligned} m_{i;\beta} &\equiv \mathbb{E} \left[z_{i;\beta}^{(L)}(T) \right] \\ &= m_{i;\beta}^{\text{NTK}} + \frac{1}{n_{L-1}} \left(m_{i;\beta}^{\Delta \text{NTK}} + m_{i;\beta}^{\text{dNTK}} + m_{i;\beta}^{\text{ddNTK-I}} + m_{i;\beta}^{\text{ddNTK-II}} \right) \\ &\quad - \frac{1}{n_{L-1}} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\beta \tilde{\alpha}_1}^{(L)} \tilde{H}_{\tilde{\alpha}_1 \tilde{\alpha}_2}^{(L)} \left(m_{i;\tilde{\alpha}_2}^{\Delta \text{NTK}} + m_{i;\tilde{\alpha}_2}^{\text{dNTK}} + m_{i;\tilde{\alpha}_2}^{\text{ddNTK-I}} + m_{i;\tilde{\alpha}_2}^{\text{ddNTK-II}} \right) \\ &\quad + O\left(\frac{1}{n^2}\right), \end{aligned} \quad (\infty.160)$$

where the first term is the (neural tangent) kernel prediction

$$m_{i;\beta}^{\text{NTK}} \equiv \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\beta \tilde{\alpha}_1}^{(L)} \tilde{H}_{\tilde{\alpha}_1 \tilde{\alpha}_2}^{(L)} y_{i;\tilde{\alpha}_2}, \quad (\infty.161)$$

and the four other kinds of terms come from the leading-order finite-width correction. Specifically, (i) the fluctuation of the NTK gives

$$m_{i;\delta}^{\Delta \text{NTK}} \equiv \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4} \left(A_{(\delta \tilde{\alpha}_1)(\tilde{\alpha}_2 \tilde{\alpha}_3)}^{(L)} + B_{\delta \tilde{\alpha}_2 \tilde{\alpha}_1 \tilde{\alpha}_3}^{(L)} + n_L B_{\delta \tilde{\alpha}_3 \tilde{\alpha}_1 \tilde{\alpha}_2}^{(L)} \right) \tilde{H}_{(L)}^{\tilde{\alpha}_1 \tilde{\alpha}_2} \tilde{H}_{(L)}^{\tilde{\alpha}_3 \tilde{\alpha}_4} y_{i;\tilde{\alpha}_4}, \quad (\infty.162)$$

where we used (8.82) to evaluate the NTK variance in terms of our decomposition into tensors $A^{(L)}$ and $B^{(L)}$; (ii) the dNTK gives

$$\begin{aligned} m_{i;\delta}^{\text{dNTK}} &\equiv - \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4} \left[2 \left(P_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + Q_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + n_L Q_{\delta \tilde{\alpha}_2 \tilde{\alpha}_1 \tilde{\alpha}_3}^{(L)} \right) Z_{\text{B}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \right. \\ &\quad + \left(n_L P_{\tilde{\alpha}_1 \delta \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + Q_{\tilde{\alpha}_1 \delta \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + Q_{\tilde{\alpha}_1 \tilde{\alpha}_2 \delta \tilde{\alpha}_3}^{(L)} \right) Z_{\text{A}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \\ &\quad \left. + \left(P_{\tilde{\alpha}_1 \delta \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + n_L Q_{\tilde{\alpha}_1 \delta \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} + Q_{\tilde{\alpha}_1 \tilde{\alpha}_2 \delta \tilde{\alpha}_3}^{(L)} \right) Z_{\text{A}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_4 \tilde{\alpha}_3} \right] y_{i;\tilde{\alpha}_4}, \end{aligned} \quad (\infty.163)$$

where we used (11.42) to evaluate the dNTK-preactivation cross correlators in terms of our decomposition into tensors $P^{(L)}$ and $Q^{(L)}$; (iii) the first ddNTK gives

$$\begin{aligned}
 m_{i;\delta}^{\text{ddNTK-I}} & \quad (\infty.164) \\
 \equiv - \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_6} & \left[R_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} \left(Z_{\text{IB}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} + Z_{\text{IB}}^{\tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_1 \tilde{\alpha}_5 \tilde{\alpha}_6 \tilde{\alpha}_4} + Z_{\text{IB}}^{\tilde{\alpha}_3 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_6 \tilde{\alpha}_4 \tilde{\alpha}_5} \right) \right. \\
 & + R_{\tilde{\alpha}_1 \delta \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} Z_{\text{IA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} + R_{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \delta}^{(L)} Z_{\text{IA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_5 \tilde{\alpha}_6 \tilde{\alpha}_4} + R_{\tilde{\alpha}_1 \tilde{\alpha}_3 \delta \tilde{\alpha}_2}^{(L)} Z_{\text{IA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_6 \tilde{\alpha}_4 \tilde{\alpha}_5} \Big] \\
 & \times \left[y_{i;\tilde{\alpha}_4} \left(\sum_j y_{j;\tilde{\alpha}_5} y_{j;\tilde{\alpha}_6} + n_L K_{\tilde{\alpha}_5 \tilde{\alpha}_6}^{(L)} \right) + y_{i;\tilde{\alpha}_5} K_{\tilde{\alpha}_6 \tilde{\alpha}_4}^{(L)} + y_{i;\tilde{\alpha}_6} K_{\tilde{\alpha}_4 \tilde{\alpha}_5}^{(L)} \right],
 \end{aligned}$$

where we used (∞.11) to evaluate the first ddNTK mean in terms of decomposition into the tensor $R^{(L)}$; and (iv) the second ddNTK gives

$$\begin{aligned}
 m_{i;\delta}^{\text{ddNTK-II}} & \quad (\infty.165) \\
 \equiv - \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_6} & \left[S_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} Z_{\text{IIB}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} + T_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} Z_{\text{IIB}}^{\tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_1 \tilde{\alpha}_5 \tilde{\alpha}_6 \tilde{\alpha}_4} + U_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3}^{(L)} Z_{\text{IIB}}^{\tilde{\alpha}_3 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_6 \tilde{\alpha}_4 \tilde{\alpha}_5} \right. \\
 & + S_{\tilde{\alpha}_1 \tilde{\alpha}_2 \delta \tilde{\alpha}_3}^{(L)} Z_{\text{IIA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_5 \tilde{\alpha}_6 \tilde{\alpha}_4} + T_{\tilde{\alpha}_1 \delta \tilde{\alpha}_3 \tilde{\alpha}_2}^{(L)} Z_{\text{IIA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4 \tilde{\alpha}_5 \tilde{\alpha}_6} + U_{\tilde{\alpha}_1 \tilde{\alpha}_3 \tilde{\alpha}_2 \delta}^{(L)} Z_{\text{IIA}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_6 \tilde{\alpha}_4 \tilde{\alpha}_5} \Big] \\
 & \times \left[y_{i;\tilde{\alpha}_4} \left(\sum_j y_{j;\tilde{\alpha}_5} y_{j;\tilde{\alpha}_6} + n_L K_{\tilde{\alpha}_5 \tilde{\alpha}_6}^{(L)} \right) + y_{i;\tilde{\alpha}_5} K_{\tilde{\alpha}_6 \tilde{\alpha}_4}^{(L)} + y_{i;\tilde{\alpha}_6} K_{\tilde{\alpha}_4 \tilde{\alpha}_5}^{(L)} \right],
 \end{aligned}$$

where we used (∞.12) to evaluate the second ddNTK mean in terms of decomposition into the tensors $S^{(L)}$, $T^{(L)}$, and $U^{(L)}$.

Interestingly, we see that not only does the mean prediction depend on the NTK differentials – as we expect, given the nontrivial representation learning at finite width – but also it depends on the NTK variance as well, cf. (∞.162). This is natural as each network fits the training data with its own *particular* NTK, and so the resulting fully-trained particular output (∞.154) depends on the NTK fluctuation, as we’ve already emphasized. In general, it is really important to understand the tradeoffs between both contributions, and in the next subsection we will return to comment on this interplay between random fluctuations and directed representation learning in the overall ensemble.

In addition, looking at the terms that are cubic in $y_{i;\tilde{\alpha}}$ in the contributions from the ddNTKs, (∞.164) and (∞.165), we see that the j -th component of the observed outputs influences the i -th component of the mean prediction for $i \neq j$. Just as we saw for the mean prediction of Bayesian inference, (6.88), this is one consequence of the wiring property of finite-width neural networks. To see another manifestation of this wiring property, let’s consider the covariance:

$$\text{Cov} \left[z_{i_1; \tilde{\beta}_1}^{(L)}(T), z_{i_2; \tilde{\beta}_2}^{(L)}(T) \right] \equiv \mathbb{E} \left[z_{i_1; \tilde{\beta}_1}^{(L)}(T) z_{i_2; \tilde{\beta}_2}^{(L)}(T) \right] - \mathbb{E} \left[z_{i_1; \tilde{\beta}_1}^{(L)}(T) \right] \mathbb{E} \left[z_{i_2; \tilde{\beta}_2}^{(L)}(T) \right]. \quad (\infty.166)$$

While we won't print this quantity in full – the full expression doesn't really play nicely with the constraints of the page – you can easily extract insight by considering some specific contributions. For instance, we can see the imprint of output-component wiring by looking at the following contribution to the covariance,

$$\begin{aligned} & \sum_{\substack{j_1, j_2 \\ \tilde{\alpha}_1, \dots, \tilde{\alpha}_4}} \mathbb{E} \left[z_{i_2; \tilde{\beta}_2}^{(L)} \widehat{\mathrm{d}H}_{i_1 j_1 j_2; \tilde{\beta}_1 \tilde{\alpha}_1 \tilde{\alpha}_2}^{(L)} \right] Z_{\mathrm{B}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \mathbb{E} \left[\left(z_{j_1; \tilde{\alpha}_3}^{(L)} - y_{j_1; \tilde{\alpha}_3} \right) \left(z_{j_2; \tilde{\alpha}_4}^{(L)} - y_{j_2; \tilde{\alpha}_4} \right) \right] \\ &= \frac{1}{n_{L-1}} \delta_{i_1 i_2} \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4} \left[\left(n_L P_{\tilde{\beta}_1 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\beta}_2}^{(L)} + Q_{\tilde{\beta}_1 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\beta}_2}^{(L)} + Q_{\tilde{\beta}_1 \tilde{\alpha}_2 \tilde{\alpha}_1 \tilde{\beta}_2}^{(L)} \right) G_{\tilde{\alpha}_3 \tilde{\alpha}_4}^{(L)} \right. \\ & \quad \left. + P_{\tilde{\beta}_1 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\beta}_2}^{(L)} \left(\sum_j y_{j; \tilde{\alpha}_3} y_{j; \tilde{\alpha}_4} \right) \right] Z_{\mathrm{B}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \\ & \quad + \frac{1}{n_{L-1}} \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4} Q_{\tilde{\beta}_1 \tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\beta}_2}^{(L)} Z_{\mathrm{B}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (y_{i_1; \tilde{\alpha}_3} y_{i_2; \tilde{\alpha}_4} + y_{i_1; \tilde{\alpha}_4} y_{i_2; \tilde{\alpha}_3}), \quad (\infty.167) \end{aligned}$$

which comes from the cross correlation between a factor of $z_{i_2; \tilde{\beta}_2}^{(L)}$ from $z_{i_2; \tilde{\beta}_2}^{(L)}(T)$ and one of the dNTK terms from $z_{i_1; \tilde{\beta}_1}^{(L)}(T)$ that involves the algorithm projector $Z_{\mathrm{B}}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4}$. In particular, we see that wiring is exhibited in the last term on the final line when the true outputs have nonzero components for both i_1 and i_2 . In this case, this correlation ($\infty.167$) implies that the test-set predictions for $z_{i_1; \tilde{\beta}_1}^{(L)}(T)$ can be shifted given $z_{i_2; \tilde{\beta}_2}^{(L)}(T)$, for $i_1 \neq i_2$: cf. our related discussion of Hebbian learning in terms of the *fire-together* inductive bias in the finite-width prior in §6.4.1 and then our following discussion of how that leads to the posterior distribution's *wiring together* in §6.4.2. Here, the presence of this wiring in the covariance of the solution means that this contribution to wiring occurs differently for each network in the ensemble: the fluctuations from instantiation to instantiation of the parameters at initialization break the permutation symmetry among the output neurons. This type of wiring suggests that each network is able to use the fluctuations between the different initial output components to its advantage in order to correlate such components together over the course of learning.

More broadly, unlike the kernel prediction at infinite width (10.39), the prediction at finite width ($\infty.154$) now has non-Gaussian statistics. In particular, since finite-width prediction is a nontrivial functional of the network output, the NTK, the dNTK, and the ddNTKs – all at initialization – and since we know that the joint distribution of those quantities $p(z^{(L)}, \widehat{H}^{(L)}, \widehat{\mathrm{d}H}^{(L)}, \widehat{\mathrm{dd}H}^{(L)}, \widehat{\mathrm{dd}H}^{(L)} | \mathcal{D})$ is a *nearly-Gaussian distribution*, then so is the fully-trained distribution $p(z^{(L)}(T) | \mathcal{D})$. In particular, there are nontrivial higher-point connected correlators; the explicit expressions of such correlators are challenging to display in any media format, though all the information needed to do so is contained in ($\infty.154$), and it's not that hard to zero in on any particular term of interest. The information carried in such correlators probably contains useful insight into some of the behavior of fully-trained networks and is likely worth further consideration.

Generalization at Finite Width

Having discussed many of the qualitative differences between the finite-width prediction (∞.154) and the infinite-width kernel prediction (10.39), now let's make some quantitative statements. In order to get a high-level understanding of how generalization is modified at finite width as compared to our extensive infinite-width analysis in §10.3, we need to determine the size of the relative importance of the finite-width corrections to our predictions, namely how these corrections depend on the widths of the hidden layers n_ℓ and the depth of the network L . In particular, we want to understand corrections to our generalized bias–variance tradeoff (10.52) at finite width.

Let's start by considering the bias term, $m_{i;\hat{\beta}} - y_{i;\hat{\beta}}$, for which we just need to look at the mean prediction $m_{i;\hat{\beta}}$ in (∞.160). In particular, we need to compare the first term, (∞.161), which more or less corresponds to the infinite-width kernel prediction $m_{i;\hat{\beta}}^\infty$ – up to the subleading and unimportant correction to the NTK mean – against the other set of finite-width contributions to the mean, (∞.162)–(∞.165).

Now, among the finite-width contributions in (∞.160), the terms inside the first set of large parentheses have a sample index corresponding to a test input, $\hat{\beta}$, in each of the tensors $A^{(L)}$, $B^{(L)}$, $P^{(L)}$, $Q^{(L)}$, $R^{(L)}$, $S^{(L)}$, $T^{(L)}$, and $U^{(L)}$. Thus, even for a training set with one training sample as we studied in §10.3.1, the particular details of these terms require an understanding of asymptotic multiple-input solutions of the recursions for the NTK variance, the preactivation–dNTK cross correlation, and the ddNTK mean; such an analysis is kind of annoying and was previously left as an adventure for thrill seekers, with a brief instruction manual for those interested buried in footnote 8 of §11.3. Unfortunately, we have no plans here to update that manual any further, and we don't expect to miss much by not doing so.²⁴

In contrast, the terms inside the second set of large parentheses in (∞.160) can be analyzed with the results we already have in hand. Thus, to nonlinearly gain a large amount of intuition with a small amount of effort, let's compare the infinite-width term (∞.161) against only this last set of terms. In particular, since both of these terms are preceded by a common prefactor $\sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} H_{\beta\tilde{\alpha}_1}^{(L)} \tilde{H}_{(L)}^{\tilde{\alpha}_1\tilde{\alpha}_2}$ – whose effect was analyzed in §10.3 – we simply need to understand the depth and width dependence of the tensors inside the second set of brackets in (∞.160). Here we'll evaluate these tensors for a single training set input $x_{\tilde{\alpha}} = x$, dropping the training sample indices for the rest of this analysis for notational simplicity.

To see the physics, it's simplest to look at the ODE limit of many many steps of gradient descent, $\eta \searrow 0$, for which the algorithm projectors take particularly simple form,

²⁴In particular, we expect that this first set of terms will be qualitatively similar to the second set of terms that we will discuss in the next paragraph, though the particular order-one-level details may differ.

$$Z_A = Z_B = \frac{1}{2 \left(\tilde{H}^{(L)} \right)^2}, \quad Z_{IA} = Z_{IB} = Z_{IIA} = \frac{-1}{6 \left(\tilde{H}^{(L)} \right)^3}, \quad Z_{IIB} = \frac{-1}{2 \left(\tilde{H}^{(L)} \right)^3}; \quad (\infty.168)$$

cf. $(\infty.138)$, $(\infty.139)$, and $(\infty.149)$ – $(\infty.153)$ to derive these expressions.²⁵ Substituting these projectors $(\infty.168)$ into the dNTK and ddNTK contributions to the mean prediction, $(\infty.163)$ – $(\infty.165)$, and then considering the $1/n$ part of the mean prediction $m_{i;\hat{\beta}}$ $(\infty.160)$, we see that all the individual terms are proportional to one of the following dimensionless ratios:

$$\begin{aligned} & \frac{A^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^2}, \quad \frac{B^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^2}, \quad \frac{P^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^2}, \quad \frac{Q^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^2}, \quad (\infty.169) \\ & \frac{R^{(L)} K^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^3}, \quad \frac{S^{(L)} K^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^3}, \quad \frac{T^{(L)} K^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^3}, \quad \frac{U^{(L)} K^{(L)}}{n_{L-1} \left(\tilde{H}^{(L)} \right)^3}. \end{aligned}$$

Thus, these eight ratios determine the size of the finite-width effects as compared to the leading infinite-width term.²⁶

Importantly, recalling our scaling laws for the NTK variance, (9.27) , for the dNTK–preactivation cross correlation, (11.65) , and for the ddNTK–preactivation cross correlations, $(\infty.21)$ and $(\infty.22)$, we see that each of these dimensionless ratios will scale like the depth-to-width ratio L/n (except for the subdominant $U^{(L)}$ contribution). Thus, overall we should find for the finite-width corrections that

$$m_{i;\hat{\beta}} - m_{i;\hat{\beta}}^\infty = O\left(\frac{L}{n}\right), \quad (\infty.170)$$

²⁵Backing off the ODE limit, for many many steps of gradient descent with a finite learning rate η , the single-input dNTK algorithm projectors, $(\infty.142)$ and $(\infty.143)$, are instead given by

$$Z_A = \frac{1}{\left(\tilde{H}^{(L)} \right)^2} \left(\frac{1 - \eta \tilde{H}^{(L)}}{2 - \eta \tilde{H}^{(L)}} \right), \quad Z_B = \frac{1}{2 \left(\tilde{H}^{(L)} \right)^2}. \quad (\infty.171)$$

In conjunction with similar single-input limits of the ddNTK algorithm projectors, this will mean that the same set of ratios, $(\infty.169)$, determine the generalization error bias term, though now with an additional η dependence. In particular, the projector Z_A diverges as the global learning rate approaches from below $\eta \nearrow 2/\tilde{H}^{(L)}$. This divergence is expected: as we noted in footnote 16, we need $\|I - \eta \hat{H}\|_\infty < 1$ for the training dynamics to converge after many many steps, $t \rightarrow \infty$. If you check, you’ll also see that some of the ddNTK algorithm projectors have the same divergence.

²⁶Incidentally, the form of the final two ratios on the first line of $(\infty.169)$ is the ultimate justification for why in (11.62) we divided the dNTK–preactivation cross-correlation tensors by two factors of the NTK mean, cf. our discussion of dimensional analysis around (11.64) . Similarly the form of the four ratios on the second line of $(\infty.169)$ is the ultimate justification for the ratios $(\infty.21)$ and $(\infty.22)$. For this latter set of ratios, we should have really written $K^{(L)} + n_L^{-1} \sum_j y_j^2$ instead of just $K^{(L)}$, especially when $K^{(L)}$ behaves as a nontrivial power law in L ; in such cases, we should really rescale the target output $y_{i;\hat{\alpha}}$ by the power-law factor $L^{-p_0/2}$ as discussed around $(\infty.26)$.

where the exact order-one constant is not important. What *is* important is that we confirmed our expectation for the overall scaling of this correction, going as the *cutoff* of our effective theory: $r \equiv L/n$. Similarly, we could do an analysis for the variance term of the generalization error by looking at the covariance (∞.166), which will be a whole lot of work with very little payoff; such an analysis would merely confirm that the leading finite-width corrections will again be of order $O(L/n)$.

In general, given that the aspect ratio L/n controls both the fluctuations in the ensemble *and* representation learning, the optimal value of the ratio is likely nonzero but also small. In particular, representation learning is enhanced by the depth, but networks with too large a value of L/n will have an effect on the mean prediction of the ensemble, but perhaps even more importantly will lead to exponentially-large problems when working with only a *particular* network: for large enough L/n , our principle of typicality can break down, and so the generalization error can begin to exhibit exponential behavior.²⁷

This further explains the success of our effective theory description at $O(L/n)$: a description with vanishing L/n , i.e., the infinite-width limit, is too simple to model the properties of deep neural networks in practice; a description for larger values of L/n , i.e., a *small*-width or *overly*-deep regime that includes many higher-order corrections, describes networks that are unlikely to be trainable; a description with small but nonzero L/n , i.e., the leading finite-width effective theory accurate to $O(L/n)$, is as simple as it can be and still accurately describe the networks that work well in practice.²⁸

In summary, we've seen that the leading finite-width corrections all scale exactly according to our long-standing expectations, and we've discussed at a high level the potential tradeoff of depth versus width. In principle, we could go further in our analysis, evaluating the multi-input recursions for nearby inputs, evaluating all the terms in the covariance, and finding all the $O(L/n)$ contributions to the generalization error with specific coefficients.²⁹ If we did this, what could it tell us? In particular, can we theoretically optimize the aspect ratio L/n for a particular activation function without any experimentation?

Unfortunately at the order we're working, we won't be able to optimize the aspect ratio L/n using the prediction formula (∞.154) from the end of training: the linear

²⁷In our discussion of *fluctuations* in §5.4, we explained that too large a value of L/n may lead to a difficult fine-tuning problem in terms of getting a particular network to behave critically.

On the one hand, the contribution of such fluctuations to the bias part of the generalization error is one way to see how the downstream effect of such fluctuations may lead to problems after training. On the other hand, the fact that fluctuations can ruin criticality for a particular network is another problem. The former problem affects the ensemble as a whole, while the latter problem affects individual networks.

²⁸To better understand the scale that separates the trainable regime from the overly-deep regime, see Appendix A. For a discussion of how to extend this trainable regime to greater depths for fixed width, see our discussion of residual networks in Appendix B.

²⁹You can also more easily evaluate the multi-input version of the relevant recursions numerically with this scaling in mind in order to determine the overall coefficient for a particular activation function of interest. This would be one way to analyze these solutions for the ReLU-like GELU and SWISH networks.

dependence of the generalization on the ratio means that its derivative is independent of the ratio; instead, we'd need to compute the higher-order corrections of order $O(L^2/n^2)$ to optimize the ratio – which is hard, though at least straightforward to do in the effective theory framework we've developed.³⁰ At order $O(L/n)$, the best we could hope to see is whether nonzero L/n improves generalization or not by looking at the sign of its overall coefficient. Rather than going through that hassle, we'll instead get a little more mileage out of the mean prediction $m_{i;\hat{\beta}}$ by trying to understand how the *algorithm dependence* at finite width leads to additional tradeoffs within the bias term of the generalization error.

Inductive Bias of the Training Algorithm

As multiply mentioned, one of the key differences between the infinite-width and finite-width networks optimized via gradient-based learning is the *algorithm dependence* of the fully-trained solution for finite-width networks. In particular, the solution – at least for MSE losses – takes a universal form, (∞.154), with all of the dependence of the solution on the training algorithm encoded in the *algorithm projectors*, Z_A , Z_B , Z_{IA} , Z_{IB} , Z_{IIA} , and Z_{IIB} , whose functional forms we've established explicitly for the two second-order updates in (∞.148), for many many steps of gradient descent in (∞.142)–(∞.147), and for the ODE limit of gradient descent in (∞.149)–(∞.151). Through this universal projection, we now have a theoretical means of isolating the *inductive bias of the training algorithm* from all the other inductive biases that are present in a fully-trained neural network. This provides a natural way of theoretically evaluating the relative merits of different training algorithms.

One aspect that is apparent just from staring at the stochastic prediction (∞.154) is that the algorithm projectors only induce projections on the NTK-differential part of the finite-width prediction and cannot affect the NTK-fluctuation contribution. More concretely, let's continue our discussion of the bias term in the generalization error. In particular, the bias is given by the following difference:

$$m_{i;\hat{\beta}} - y_{i;\hat{\beta}}. \quad (\infty.172)$$

On the one hand, it's clear that the NTK-variance contribution to the mean prediction, (∞.162), encoded in $A^{(L)}$ and $B^{(L)}$ is irreducible, entirely independent of the training algorithm; this irreducible NTK-variance contribution depends on the training data in a fixed way and arises due to instantiation-to-instantiation fluctuations in the NTK across different realizations of the parameters. On the other hand, the projectors always act on

³⁰In §A.3, we'll compute such higher-order effects from an *information-theoretic* perspective. This analysis will give a heuristic prescription for optimizing the network's depth-to-width ratio r : we consider an auxiliary *unsupervised learning* objective thought to be beneficial for building representations and optimize the aspect ratio in terms of *that* criterion rather than the generalization error. In this setting, we're able to determine optimal aspect ratios for different activation functions by trading off leading-order effects against higher-order effects. This is somewhat similar in spirit to the way in which Newton's method trades off the leading decrease of the loss against the subleading loss increase in order to optimize the overall learning rate, cf. our discussion of Newton's method in footnote 8 of §10.2.

the NTK-differential tensors $P^{(L)}$, $Q^{(L)}$, $R^{(L)}$, $S^{(L)}$, $T^{(L)}$, and $U^{(L)}$; this means that the dNTK's and ddNTKs' effect on the network's predictions is adjustable by the training algorithm and can be tuned differently for different datasets and tasks.

To reiterate our previous discussion of the tradeoff, on the one hand the NTK-fluctuation contribution is likely harmful as it leads to a breakdown of criticality for any *particular* network. On the other hand, the latter NTK-differential contribution is the ultimate source of the nontrivial representation learning at finite width, and so it would be nice to make this contribution large. With these algorithm projectors, we now have direct means to enhance the latter benefit while keeping fixed the former cost.

One way of understanding these algorithm projectors is as the sample-space *dual description* of a learning algorithm, analogous to the relationship between the feature functions $\phi_j(x)$ and the kernel $k(x_{\delta_1}, x_{\delta_2})$ that we explored in §10.4 and §11.4. From the parameter-space microscopic perspective, a learning algorithm, such as gradient descent, explicitly operates on the parameters. At the end of training, all the details of the algorithm are left implicit in the trained parameters θ^* , making it difficult to understand its effect on the model's predictions. From this perspective, the algorithm is simple to define (§7.2) but decidedly difficult to analyze for general interacting machine-learning models (§ $\infty.2.2$). However, if you can analyze or *solve* a model to find its sample-space dual description, then the algorithm projectors make the influence of the learning algorithm explicit.³¹

Analogously to the (nearly-)kernel perspective on (nearly-)linear models, the appearance of the algorithm projectors in the generalization error means that we can begin to think about *engineering* them directly, either by picking them to be directly used with a prediction formula such as ($\infty.154$) or alternatively by attempting to find the parameter-space dual of an algorithm projector. While we expect this latter task to be difficult – just as it's often hard to find the feature functions that correspond to a particular kernel – it could be very worthwhile to explore as a means of engineering the inductive bias of the training algorithm directly.³²

In particular, we expect that this could significantly improve generalization by designing algorithms that are better tailored to the details of an architecture or the properties of the underlying dataset or task.³³ This design process could in principle proceed as follows: (i) engineer a desired functional form for the algorithm projectors and then (ii) determine the parameter-space training algorithm that leads to the projectors having those desired properties or specific functional form. While further details of such **inverse algorithm design** is outside the scope of this book, we think that this line of *dual* analysis has the potential to unlock a much deeper understanding of the relationships

³¹This is the sense in which we meant that the conditioning on the learning algorithm would be *simple* for the trained ensemble in (0.9).

³²This can also be seen as another advantage of nonlinear or interacting machine learning models. For linear models, the solution is always independent of the details of the learning algorithm (§10.2.2).

³³Since the algorithm projectors have sample indices, the *optimal* choice of the algorithm will in general depend on the details and structure of the training set in addition to the definition of the model.

between the inductive bias of the training algorithm, the inductive bias of the network architecture, and the ultimate success of the model.

There's No Place Where Gradient Descent = Exact Bayesian Inference

Finally, a curious reader might wonder whether the connection that we detailed in §10.2.4 between gradient descent optimization and exact Bayesian inference for the infinite-width limit persists for more realistic networks at finite width. In short, the answer is no.

In long, recall the training hyperparameter settings (10.46) and (10.48) that matched gradient-based learning with Bayesian inference at infinite width. In particular, the latter condition (10.48) required turning off any learning in the hidden layers:

$$\lambda_b^{(\ell)} = 0, \quad \lambda_W^{(\ell)} = 0, \quad \text{for } \ell < L. \quad (\infty.173)$$

Importantly, this makes the hidden-layer NTK vanish exactly: $\hat{H}_{i_1 i_2; \delta_1 \delta_2}^{(\ell)} = 0$, for $\ell < L$ at any finite width, cf. (8.12). Next, flip back and look at the P - and Q -recursions for the dNTK–preactivation cross-correlation tensors, (11.49) and (11.52), and then you'll probably also want to flip further back and visit the F - and B -recursions, (8.79) and (8.89), respectively. (We'll be waiting here for you when you return.) Immediately, you should notice a problem: if the hidden-layer NTK mean vanishes, then at this order $B^{(\ell)} = 0$, $F^{(\ell)} = 0$, and all this together implies that $P^{(L)} = Q^{(L)} = 0$. Thus, all the effects of the dNTK are turned off. Similarly, if you flip forth and look at the R -, S -, T -, and U -recursions, ($\infty.177$) and ($\infty.179$)–($\infty.181$), then you'll notice that all the effects of the ddNTKs are turned off as well. Thus, the mean prediction of our finite-width gradient-based-learning ensemble ($\infty.160$) *cannot* match the exact Bayesian posterior mean at finite width (6.88).

Note that this mismatch is obvious in hindsight; by only training the last layer ($\infty.173$), we get a linear model (10.134). In other words, since the hyperparameter choices of ($\infty.173$) lead to a model with random features that are fixed over the course of training, there's no representation learning possible for these settings. In contrast, we found nontrivial representation learning when studying exact Bayesian inference at finite width in §6.4.3.

Ultimately, for Bayesian inference we only care about the preactivation distribution $p(z^{(\ell)}|\mathcal{D})$, while for gradient-based learning we need to consider the joint preactivation–NTK–dNTK–ddNTKs distribution $p(z^{(L)}, \hat{H}^{(L)}, \widehat{dH}^{(L)}, \widehat{dd_I H}^{(L)}, \widehat{dd_{II} H}^{(L)}|\mathcal{D})$, which incorporates the statistics of the derivatives of the preactivations at initialization: such derivatives of the model output are invisible to the exact Bayesian inferencer.

$\infty.3$ RG Flow of the ddNTKs: The Full Expressions

These expressions were kind of horrible, so we decided to hide them here at the end of the chapter. As such, they are only really needed for three reasons: (i) to explicitly check the absence of any NTK differentials for the hyperparameter setup ($\infty.173$) and thus confirm

that the connection between gradient descent and Bayesian inference doesn't persist at finite width, (ii) to check the details of the depth-to-width scaling of the ddNTKs that we discussed in §∞.1, and (iii) to more generally evaluate the ddNTKs' contributions to the ensemble's mean prediction, (∞.164) and (∞.165), for multiple inputs – analytically or numerically – or to compute other higher-order statistics of our stochastic prediction (∞.154).

$\widehat{\text{dd}_I H}$ Stochastic Forward Equation

$$\begin{aligned} & \widehat{\text{dd}_I H}_{i_0 i_1 i_2 i_3; \delta_0 \delta_1 \delta_2 \delta_3}^{(\ell+1)} \\ &= \delta_{i_0 i_1} \frac{\lambda_W^{(\ell+1)}}{n_\ell} \sum_{j_0, j_1, j_2, j_3=1}^{n_\ell} \delta_{j_0 j_1} W_{i_2 j_2}^{(\ell+1)} W_{i_3 j_3}^{(\ell+1)} \sigma_{j_1; \delta_1}^{(\ell)} \sigma_{j_2; \delta_2}'^{(\ell)} \sigma_{j_3; \delta_3}'^{(\ell)} \\ & \quad \times \left[\sigma_{j_0; \delta_0}''^{(\ell)} \widehat{H}_{j_0 j_2; \delta_0 \delta_2}^{(\ell)} \widehat{H}_{j_0 j_3; \delta_0 \delta_3}^{(\ell)} + \sigma_{j_0; \delta_0}'^{(\ell)} \widehat{\text{d}H}_{j_0 j_2 j_3; \delta_0 \delta_2 \delta_3}^{(\ell)} \right] \\ &+ \delta_{i_0 i_2} \frac{\lambda_W^{(\ell+1)}}{n_\ell} \sum_{j_0, j_1, j_2, j_3=1}^{n_\ell} \delta_{j_0 j_2} W_{i_3 j_3}^{(\ell+1)} W_{i_1 j_1}^{(\ell+1)} \sigma_{j_2; \delta_2}^{(\ell)} \sigma_{j_3; \delta_3}'^{(\ell)} \sigma_{j_1; \delta_1}'^{(\ell)} \\ & \quad \times \left[\sigma_{j_0; \delta_0}''^{(\ell)} \widehat{H}_{j_0 j_3; \delta_0 \delta_3}^{(\ell)} \widehat{H}_{j_0 j_1; \delta_0 \delta_1}^{(\ell)} + \sigma_{j_0; \delta_0}'^{(\ell)} \widehat{\text{d}H}_{j_0 j_3 j_1; \delta_0 \delta_3 \delta_1}^{(\ell)} \right] \\ &+ \delta_{i_0 i_3} \frac{\lambda_W^{(\ell+1)}}{n_\ell} \sum_{j_0, j_1, j_2, j_3=1}^{n_\ell} \delta_{j_0 j_3} W_{i_1 j_1}^{(\ell+1)} W_{i_2 j_2}^{(\ell+1)} \sigma_{j_3; \delta_3}^{(\ell)} \sigma_{j_1; \delta_1}'^{(\ell)} \sigma_{j_2; \delta_2}'^{(\ell)} \\ & \quad \times \left[\sigma_{j_0; \delta_0}''^{(\ell)} \widehat{H}_{j_0 j_1; \delta_0 \delta_1}^{(\ell)} \widehat{H}_{j_0 j_2; \delta_0 \delta_2}^{(\ell)} + \sigma_{j_0; \delta_0}'^{(\ell)} \widehat{\text{d}H}_{j_0 j_1 j_2; \delta_0 \delta_1 \delta_2}^{(\ell)} \right] \\ &+ \sum_{j_0, j_1, j_2, j_3=1}^{n_\ell} W_{i_0 j_0}^{(\ell+1)} W_{i_1 j_1}^{(\ell+1)} W_{i_2 j_2}^{(\ell+1)} W_{i_3 j_3}^{(\ell+1)} \sigma_{j_1; \delta_1}'^{(\ell)} \sigma_{j_2; \delta_2}'^{(\ell)} \sigma_{j_3; \delta_3}'^{(\ell)} \\ & \quad \times \left[\sigma_{j_0; \delta_0}'^{(\ell)} \widehat{\text{dd}_I H}_{j_0 j_1 j_2 j_3; \delta_0 \delta_1 \delta_2 \delta_3}^{(\ell)} + \sigma_{j_0; \delta_0}''^{(\ell)} \widehat{\text{d}H}_{j_0 j_1 j_2; \delta_0 \delta_1 \delta_2}^{(\ell)} \widehat{H}_{j_0 j_3; \delta_0 \delta_3}^{(\ell)} \right. \\ & \quad + \sigma_{j_0; \delta_0}''^{(\ell)} \widehat{\text{d}H}_{j_0 j_2 j_3; \delta_0 \delta_2 \delta_3}^{(\ell)} \widehat{H}_{j_0 j_1; \delta_0 \delta_1}^{(\ell)} + \sigma_{j_0; \delta_0}''^{(\ell)} \widehat{\text{d}H}_{j_0 j_3 j_1; \delta_0 \delta_3 \delta_1}^{(\ell)} \widehat{H}_{j_0 j_2; \delta_0 \delta_2}^{(\ell)} \\ & \quad \left. + \sigma_{j_0; \delta_0}'''^{(\ell)} \widehat{H}_{j_0 j_1; \delta_0 \delta_1}^{(\ell)} \widehat{H}_{j_0 j_2; \delta_0 \delta_2}^{(\ell)} \widehat{H}_{j_0 j_3; \delta_0 \delta_3}^{(\ell)} \right]. \end{aligned} \tag{\infty.174}$$

$\widehat{\text{dd}_{II} H}$ Stochastic Forward Equation

$$\begin{aligned} & \widehat{\text{dd}_{II} H}_{i_1 i_2 i_3 i_4; \delta_1 \delta_2 \delta_3 \delta_4}^{(\ell+1)} \\ &= \delta_{i_1 i_3} \delta_{i_2 i_4} \left(\frac{\lambda_W^{(\ell+1)}}{n_\ell} \right)^2 \sum_{j, k=1}^{n_\ell} \sigma_{j; \delta_1}'^{(\ell)} \sigma_{k; \delta_2}'^{(\ell)} \sigma_{j; \delta_3}^{(\ell)} \sigma_{k; \delta_4}^{(\ell)} \widehat{H}_{j k; \delta_1 \delta_2}^{(\ell)} \\ &+ \delta_{i_1 i_2} \frac{\lambda_W^{(\ell+1)}}{n_\ell} \sum_{j_1, \dots, j_4=1}^{n_\ell} \delta_{j_1 j_2} W_{i_3 j_3}^{(\ell+1)} W_{i_4 j_4}^{(\ell+1)} \sigma_{j_1; \delta_1}'^{(\ell)} \sigma_{j_2; \delta_2}'^{(\ell)} \sigma_{j_3; \delta_3}'^{(\ell)} \sigma_{j_4; \delta_4}'^{(\ell)} \widehat{H}_{j_1 j_3; \delta_1 \delta_3}^{(\ell)} \widehat{H}_{j_2 j_4; \delta_2 \delta_4}^{(\ell)} \end{aligned} \tag{\infty.175}$$

$$\begin{aligned}
 & + \delta_{i_1 i_3} \frac{\lambda_W^{(\ell+1)}}{n_\ell} \sum_{j_1, \dots, j_4=1}^{n_\ell} \delta_{j_1 j_3} W_{i_2 j_2}^{(\ell+1)} W_{i_4 j_4}^{(\ell+1)} \sigma_{j_3; \delta_3}^{(\ell)} \sigma_{j_4; \delta_4}'^{(\ell)} \sigma_{j_1; \delta_1}'^{(\ell)} \\
 & \quad \times \left[\sigma_{j_2; \delta_2}''^{(\ell)} \widehat{H}_{j_2 j_1; \delta_2 \delta_1}^{(\ell)} \widehat{H}_{j_2 j_4; \delta_2 \delta_4}^{(\ell)} + \sigma_{j_2; \delta_2}'^{(\ell)} \widehat{\text{d}H}_{j_2 j_1 j_4; \delta_2 \delta_1 \delta_4}^{(\ell)} \right] \\
 & + \delta_{i_2 i_4} \frac{\lambda_W^{(\ell+1)}}{n_\ell} \sum_{j_1, \dots, j_4=1}^{n_\ell} \delta_{j_2 j_4} W_{i_1 j_1}^{(\ell+1)} W_{i_3 j_3}^{(\ell+1)} \sigma_{j_4; \delta_4}^{(\ell)} \sigma_{j_3; \delta_3}'^{(\ell)} \sigma_{j_2; \delta_2}'^{(\ell)} \\
 & \quad \times \left[\sigma_{j_1; \delta_1}''^{(\ell)} \widehat{H}_{j_1 j_2; \delta_1 \delta_2}^{(\ell)} \widehat{H}_{j_1 j_3; \delta_1 \delta_3}^{(\ell)} + \sigma_{j_1; \delta_1}'^{(\ell)} \widehat{\text{d}H}_{j_1 j_2 j_3; \delta_1 \delta_2 \delta_3}^{(\ell)} \right] \\
 & + \sum_{j_1, j_2, j_3, j_4=1}^{n_\ell} W_{i_1 j_1}^{(\ell+1)} W_{i_2 j_2}^{(\ell+1)} W_{i_3 j_3}^{(\ell+1)} W_{i_4 j_4}^{(\ell+1)} \sigma_{j_3; \delta_3}'^{(\ell)} \sigma_{j_4; \delta_4}'^{(\ell)} \\
 & \quad \times \left[\sigma_{j_1; \delta_1}'^{(\ell)} \sigma_{j_2; \delta_2}'^{(\ell)} \widehat{\text{d}d_{\Pi}H}_{j_1 j_2 j_3 j_4; \delta_1 \delta_2 \delta_3 \delta_4}^{(\ell)} + \sigma_{j_1; \delta_1}''^{(\ell)} \sigma_{j_2; \delta_2}''^{(\ell)} \widehat{H}_{j_1 j_2; \delta_1 \delta_2}^{(\ell)} \widehat{H}_{j_1 j_3; \delta_1 \delta_3}^{(\ell)} \widehat{H}_{j_2 j_4; \delta_2 \delta_4}^{(\ell)} \right. \\
 & \quad \left. + \sigma_{j_1; \delta_1}'^{(\ell)} \sigma_{j_2; \delta_2}''^{(\ell)} \widehat{H}_{j_2 j_4; \delta_2 \delta_4}^{(\ell)} \widehat{\text{d}H}_{j_1 j_2 j_3; \delta_1 \delta_2 \delta_3}^{(\ell)} + \sigma_{j_2; \delta_2}'^{(\ell)} \sigma_{j_1; \delta_1}''^{(\ell)} \widehat{H}_{j_1 j_3; \delta_1 \delta_3}^{(\ell)} \widehat{\text{d}H}_{j_2 j_1 j_4; \delta_2 \delta_1 \delta_4}^{(\ell)} \right].
 \end{aligned}$$

$\widehat{\text{dd}_I H}$ Recursion

The mean of the first ddNTK decomposes as

$$\begin{aligned}
 & \mathbb{E} \left[\widehat{\text{dd}_I H}_{i_0 i_1 i_2 i_3; \delta_0 \delta_1 \delta_2 \delta_3}^{(\ell)} \right] \\
 & = \frac{1}{n_{\ell-1}} \left[\delta_{i_0 i_1} \delta_{i_2 i_3} R_{\delta_0 \delta_1 \delta_2 \delta_3}^{(\ell)} + \delta_{i_0 i_2} \delta_{i_3 i_1} R_{\delta_0 \delta_2 \delta_3 \delta_1}^{(\ell)} + \delta_{i_0 i_3} \delta_{i_1 i_2} R_{\delta_0 \delta_3 \delta_1 \delta_2}^{(\ell)} \right].
 \end{aligned} \tag{\infty.176}$$

The tensor $R^{(\ell)}$ satisfies the following layer-to-layer recursion:

$$\begin{aligned}
 & R_{\delta_0 \delta_1 \delta_2 \delta_3}^{(\ell+1)} \\
 & = \lambda_W^{(\ell+1)} C_W^{(\ell+1)} \langle \sigma_{\delta_0}'' \sigma_{\delta_1}' \sigma_{\delta_2}' \sigma_{\delta_3}' \rangle_{G^{(\ell)}} H_{\delta_0 \delta_2}^{(\ell)} H_{\delta_0 \delta_3}^{(\ell)} \\
 & \quad + \left(\frac{n_\ell}{n_{\ell-1}} \right) \left(C_W^{(\ell+1)} \langle \sigma_{\delta_2}' \sigma_{\delta_3}' \rangle_{G^{(\ell)}} \right) \left(\lambda_W^{(\ell+1)} \langle \sigma_{\delta_0}'' \sigma_{\delta_1}' \rangle_{G^{(\ell)}} \right) B_{\delta_0 \delta_0 \delta_2 \delta_3}^{(\ell)} \\
 & \quad + \left(\frac{n_\ell}{n_{\ell-1}} \right) \left(C_W^{(\ell+1)} \langle \sigma_{\delta_2}' \sigma_{\delta_3}' \rangle_{G^{(\ell)}} \right) \left[\lambda_W^{(\ell+1)} \left(\langle \sigma_{\delta_0}'' \sigma_{\delta_1}' \rangle_{G^{(\ell)}} P_{\delta_0 \delta_2 \delta_3 \delta_0}^{(\ell)} + \langle \sigma_{\delta_0}' \sigma_{\delta_1}' \rangle_{G^{(\ell)}} P_{\delta_0 \delta_2 \delta_3 \delta_1}^{(\ell)} \right) \right] \\
 & \quad + \left(C_W^{(\ell+1)} \right)^2 \langle \sigma_{\delta_0}''' \sigma_{\delta_1}' \sigma_{\delta_2}' \sigma_{\delta_3}' \rangle_{G^{(\ell)}} H_{\delta_0 \delta_1}^{(\ell)} H_{\delta_0 \delta_2}^{(\ell)} H_{\delta_0 \delta_3}^{(\ell)} \\
 & \quad + \left(\frac{n_\ell}{n_{\ell-1}} \right) \left(C_W^{(\ell+1)} \langle \sigma_{\delta_2}' \sigma_{\delta_3}' \rangle_{G^{(\ell)}} \right) \left(C_W^{(\ell+1)} \langle \sigma_{\delta_0}''' \sigma_{\delta_1}' \rangle_{G^{(\ell)}} \right) B_{\delta_0 \delta_0 \delta_2 \delta_3}^{(\ell)} H_{\delta_0 \delta_1}^{(\ell)} \\
 & \quad + \left(\frac{n_\ell}{n_{\ell-1}} \right) \left(C_W^{(\ell+1)} \langle \sigma_{\delta_2}' \sigma_{\delta_3}' \rangle_{G^{(\ell)}} \right) \\
 & \quad \quad \left[C_W^{(\ell+1)} \left(\langle \sigma_{\delta_0}''' \sigma_{\delta_1}' \rangle_{G^{(\ell)}} P_{\delta_0 \delta_2 \delta_3 \delta_0}^{(\ell)} + \langle \sigma_{\delta_0}'' \sigma_{\delta_1}'' \rangle_{G^{(\ell)}} P_{\delta_0 \delta_2 \delta_3 \delta_1}^{(\ell)} \right) \right] H_{\delta_0 \delta_1}^{(\ell)} \\
 & \quad + \left(\frac{n_\ell}{n_{\ell-1}} \right) \left(C_W^{(\ell+1)} \langle \sigma_{\delta_2}' \sigma_{\delta_3}' \rangle_{G^{(\ell)}} \right) \left(C_W^{(\ell+1)} \langle \sigma_{\delta_0}' \sigma_{\delta_1}' \rangle_{G^{(\ell)}} \right) R_{\delta_0 \delta_1 \delta_2 \delta_3}^{(\ell)} + O\left(\frac{1}{n}\right).
 \end{aligned} \tag{\infty.177}$$

$\widehat{\text{dd}}_{\text{II}}H$ Recursions

The mean of the second ddNTK decomposes as

$$\begin{aligned} \mathbb{E} \left[\widehat{\text{dd}}_{\text{II}}H_{i_1 i_2 i_3 i_4; \delta_1 \delta_2 \delta_3 \delta_4}^{(\ell)} \right] \\ = \frac{1}{n_{\ell-1}} \left[\delta_{i_1 i_2} \delta_{i_3 i_4} S_{\delta_1 \delta_2 \delta_3 \delta_4}^{(\ell)} + \delta_{i_1 i_3} \delta_{i_4 i_2} T_{\delta_1 \delta_3 \delta_4 \delta_2}^{(\ell)} + \delta_{i_1 i_4} \delta_{i_2 i_3} U_{\delta_1 \delta_4 \delta_2 \delta_3}^{(\ell)} \right]. \end{aligned} \quad (\infty.178)$$

The tensor $S^{(\ell)}$ satisfies the following layer-to-layer recursion:

$$\begin{aligned} S_{\delta_1 \delta_2 \delta_3 \delta_4}^{(\ell+1)} \\ = C_W^{(\ell+1)} \lambda_W^{(\ell+1)} \langle \sigma'_{\delta_1} \sigma'_{\delta_2} \sigma'_{\delta_3} \sigma'_{\delta_4} \rangle_{G^{(\ell)}} H_{\delta_1 \delta_3}^{(\ell)} H_{\delta_2 \delta_4}^{(\ell)} \\ + \left(\frac{n_{\ell}}{n_{\ell-1}} \right) \left(C_W^{(\ell+1)} \langle \sigma'_{\delta_3} \sigma'_{\delta_4} \rangle_{G^{(\ell)}} \right) \\ \left[\lambda_W^{(\ell+1)} \langle \sigma'_{\delta_1} \sigma'_{\delta_2} \rangle_{G^{(\ell)}} + C_W^{(\ell+1)} H_{\delta_1 \delta_2}^{(\ell)} \langle \sigma''_{\delta_1} \sigma''_{\delta_2} \rangle_{G^{(\ell)}} \right] B_{\delta_1 \delta_2 \delta_3 \delta_4}^{(\ell)} \\ + \left(C_W^{(\ell+1)} \right)^2 \langle \sigma''_{\delta_1} \sigma''_{\delta_2} \sigma'_{\delta_3} \sigma'_{\delta_4} \rangle_{G^{(\ell)}} H_{\delta_1 \delta_2}^{(\ell)} H_{\delta_1 \delta_3}^{(\ell)} H_{\delta_2 \delta_4}^{(\ell)} \\ + \left(\frac{n_{\ell}}{n_{\ell-1}} \right) \left(C_W^{(\ell+1)} \langle \sigma'_{\delta_1} \sigma'_{\delta_2} \rangle_{G^{(\ell)}} \right) \left(C_W^{(\ell+1)} \langle \sigma'_{\delta_3} \sigma'_{\delta_4} \rangle_{G^{(\ell)}} \right) S_{\delta_1 \delta_2 \delta_3 \delta_4}^{(\ell)} + O\left(\frac{1}{n}\right). \end{aligned} \quad (\infty.179)$$

The tensor $T^{(\ell)}$ satisfies the following layer-to-layer recursion:

$$\begin{aligned} T_{\delta_1 \delta_3 \delta_4 \delta_2}^{(\ell+1)} \\ = \left(\lambda_W^{(\ell+1)} \right)^2 \langle \sigma'_{\delta_1} \sigma'_{\delta_2} \sigma_{\delta_3} \sigma_{\delta_4} \rangle_{G^{(\ell)}} H_{\delta_1 \delta_2}^{(\ell)} \\ + \left(\frac{n_{\ell}}{n_{\ell-1}} \right) \left(\lambda_W^{(\ell+1)} \right)^2 \sum_{\delta_5, \dots, \delta_8 \in \mathcal{D}} \langle z_{\delta_5} \sigma'_{\delta_1} \sigma_{\delta_3} \rangle_{G^{(\ell)}} \langle z_{\delta_6} \sigma'_{\delta_2} \sigma_{\delta_4} \rangle_{G^{(\ell)}} G_{(\ell)}^{\delta_5 \delta_7} G_{(\ell)}^{\delta_6 \delta_8} F_{\delta_7 \delta_1 \delta_8 \delta_2}^{(\ell)} \\ + C_W^{(\ell+1)} \lambda_W^{(\ell+1)} \langle \sigma'_{\delta_1} \sigma''_{\delta_2} \sigma_{\delta_3} \sigma'_{\delta_4} \rangle_{G^{(\ell)}} H_{\delta_2 \delta_1}^{(\ell)} H_{\delta_2 \delta_4}^{(\ell)} \\ + \left(\frac{n_{\ell}}{n_{\ell-1}} \right) C_W^{(\ell+1)} \lambda_W^{(\ell+1)} H_{\delta_2 \delta_4}^{(\ell)} \sum_{\delta_5, \dots, \delta_8 \in \mathcal{D}} \langle z_{\delta_5} \sigma'_{\delta_1} \sigma_{\delta_3} \rangle_{G^{(\ell)}} \langle z_{\delta_6} \sigma''_{\delta_2} \sigma'_{\delta_4} \rangle_{G^{(\ell)}} G_{(\ell)}^{\delta_5 \delta_7} G_{(\ell)}^{\delta_6 \delta_8} F_{\delta_7 \delta_1 \delta_8 \delta_2}^{(\ell)} \\ + \left(\frac{n_{\ell}}{n_{\ell-1}} \right) C_W^{(\ell+1)} \lambda_W^{(\ell+1)} \langle \sigma'_{\delta_2} \sigma'_{\delta_4} \rangle_{G^{(\ell)}} \left(\langle \sigma''_{\delta_1} \sigma_{\delta_3} \rangle_{G^{(\ell)}} Q_{\delta_2 \delta_4 \delta_1 \delta_1}^{(\ell)} + \langle \sigma'_{\delta_1} \sigma'_{\delta_3} \rangle_{G^{(\ell)}} Q_{\delta_2 \delta_4 \delta_1 \delta_3}^{(\ell)} \right) \\ + C_W^{(\ell+1)} \lambda_W^{(\ell+1)} \langle \sigma'_{\delta_2} \sigma''_{\delta_1} \sigma_{\delta_4} \sigma'_{\delta_3} \rangle_{G^{(\ell)}} H_{\delta_1 \delta_2}^{(\ell)} H_{\delta_1 \delta_3}^{(\ell)} \\ + \left(\frac{n_{\ell}}{n_{\ell-1}} \right) C_W^{(\ell+1)} \lambda_W^{(\ell+1)} H_{\delta_1 \delta_3}^{(\ell)} \sum_{\delta_5, \dots, \delta_8 \in \mathcal{D}} \langle z_{\delta_5} \sigma'_{\delta_2} \sigma_{\delta_4} \rangle_{G^{(\ell)}} \langle z_{\delta_6} \sigma''_{\delta_1} \sigma'_{\delta_3} \rangle_{G^{(\ell)}} G_{(\ell)}^{\delta_5 \delta_7} G_{(\ell)}^{\delta_6 \delta_8} F_{\delta_7 \delta_2 \delta_8 \delta_1}^{(\ell)} \\ + \left(\frac{n_{\ell}}{n_{\ell-1}} \right) C_W^{(\ell+1)} \lambda_W^{(\ell+1)} \langle \sigma'_{\delta_1} \sigma'_{\delta_3} \rangle_{G^{(\ell)}} \left(\langle \sigma''_{\delta_2} \sigma_{\delta_4} \rangle_{G^{(\ell)}} Q_{\delta_1 \delta_3 \delta_2 \delta_2}^{(\ell)} + \langle \sigma'_{\delta_2} \sigma'_{\delta_4} \rangle_{G^{(\ell)}} Q_{\delta_1 \delta_3 \delta_2 \delta_4}^{(\ell)} \right) \end{aligned} \quad (\infty.180)$$

$$\begin{aligned}
& + \left(\frac{n_\ell}{n_{\ell-1}} \right) \left(C_W^{(\ell+1)} \right)^2 \langle \sigma'_{\delta_1} \sigma'_{\delta_3} \rangle_{G^{(\ell)}} \langle \sigma'_{\delta_2} \sigma'_{\delta_4} \rangle_{G^{(\ell)}} T_{\delta_1 \delta_3 \delta_4 \delta_2}^{(\ell)} \\
& + \left(C_W^{(\ell+1)} \right)^2 \langle \sigma''_{\delta_1} \sigma''_{\delta_2} \sigma'_{\delta_3} \sigma'_{\delta_4} \rangle_{G^{(\ell)}} H_{\delta_1 \delta_2}^{(\ell)} H_{\delta_1 \delta_3}^{(\ell)} H_{\delta_2 \delta_4}^{(\ell)} \\
& + \left(\frac{n_\ell}{n_{\ell-1}} \right) \left(C_W^{(\ell+1)} \right)^2 H_{\delta_1 \delta_3}^{(\ell)} H_{\delta_2 \delta_4}^{(\ell)} \\
& \quad \sum_{\delta_5, \dots, \delta_8 \in \mathcal{D}} \langle z_{\delta_5} \sigma''_{\delta_1} \sigma'_{\delta_3} \rangle_{G^{(\ell)}} \langle z_{\delta_6} \sigma''_{\delta_2} \sigma'_{\delta_4} \rangle_{G^{(\ell)}} G_{(\ell)}^{\delta_5 \delta_7} G_{(\ell)}^{\delta_6 \delta_8} F_{\delta_7 \delta_1 \delta_8 \delta_2}^{(\ell)} \\
& + \left(\frac{n_\ell}{n_{\ell-1}} \right) \left(C_W^{(\ell+1)} \right)^2 H_{\delta_2 \delta_4}^{(\ell)} \langle \sigma'_{\delta_1} \sigma'_{\delta_3} \rangle_{G^{(\ell)}} \left(\langle \sigma'''_{\delta_2} \sigma'_{\delta_4} \rangle_{G^{(\ell)}} Q_{\delta_1 \delta_3 \delta_2 \delta_2}^{(\ell)} + \langle \sigma''_{\delta_2} \sigma''_{\delta_4} \rangle_{G^{(\ell)}} Q_{\delta_1 \delta_3 \delta_2 \delta_4}^{(\ell)} \right) \\
& + \left(\frac{n_\ell}{n_{\ell-1}} \right) \left(C_W^{(\ell+1)} \right)^2 H_{\delta_1 \delta_3}^{(\ell)} \langle \sigma'_{\delta_2} \sigma'_{\delta_4} \rangle_{G^{(\ell)}} \left(\langle \sigma'''_{\delta_1} \sigma'_{\delta_3} \rangle_{G^{(\ell)}} Q_{\delta_2 \delta_4 \delta_1 \delta_1}^{(\ell)} + \langle \sigma''_{\delta_1} \sigma''_{\delta_3} \rangle_{G^{(\ell)}} Q_{\delta_2 \delta_4 \delta_1 \delta_3}^{(\ell)} \right) \\
& + O\left(\frac{1}{n}\right).
\end{aligned}$$

The tensor $U^{(\ell)}$ satisfies the following layer-to-layer recursion:

$$\begin{aligned}
U_{\delta_1 \delta_4 \delta_2 \delta_3}^{(\ell+1)} &= \left(C_W^{(\ell+1)} \right)^2 \langle \sigma''_{\delta_1} \sigma''_{\delta_2} \sigma'_{\delta_3} \sigma'_{\delta_4} \rangle_{G^{(\ell)}} H_{\delta_1 \delta_2}^{(\ell)} H_{\delta_1 \delta_3}^{(\ell)} H_{\delta_2 \delta_4}^{(\ell)} \\
&+ \left(\frac{n_\ell}{n_{\ell-1}} \right) \left(C_W^{(\ell+1)} \langle \sigma'_{\delta_1} \sigma'_{\delta_4} \rangle_{G^{(\ell)}} \right) \left(C_W^{(\ell+1)} \langle \sigma'_{\delta_2} \sigma'_{\delta_3} \rangle_{G^{(\ell)}} \right) U_{\delta_1 \delta_4 \delta_2 \delta_3}^{(\ell)} + O\left(\frac{1}{n}\right).
\end{aligned} \tag{\infty.181}$$