
Feature Learning in Infinite-Depth Neural Networks

Greg Yang*
xAI

Dingli Yu*
Princeton Language
and Intelligence

Chen Zhu
Nvidia

Soufiane Hayou†
Simons Institute
UC Berkeley

Abstract

By classifying infinite-width neural networks and identifying the *optimal* limit, [23, 25] demonstrated a universal way, called μP , for *widthwise hyperparameter transfer*, i.e., predicting optimal hyperparameters of wide neural networks from narrow ones. Here we investigate the analogous classification for *depthwise parametrizations* of deep residual networks (resnets). We classify depthwise parametrizations of block multiplier and learning rate by their infinite-width-then-depth limits. In resnets where each block has only one layer, we identify a unique optimal parametrization, called Depth- μP that extends μP and show empirically it admits depthwise hyperparameter transfer. We identify *feature diversity* as a crucial factor in deep networks, and Depth- μP can be characterized as maximizing both feature learning and feature diversity. Exploiting this, we find that absolute value, among all homogeneous nonlinearities, maximizes feature diversity and indeed empirically leads to significantly better performance. However, if each block is deeper (such as modern transformers), then we find fundamental limitations in all possible infinite-depth limits of such parametrizations, which we illustrate both theoretically and empirically on simple networks as well as Megatron transformer trained on Common Crawl.

1 Introduction

Deep neural networks have showcased remarkable performance across a broad range of tasks, including image classification, game playing exemplified by AlphaGo [17], and natural language processing demonstrated by GPT-4 [15]. A prevailing trend in developing these networks is to increase their size and complexity, with empirical evidence indicating that using the same computation resources, models with more parameters tend to exhibit better performance. There are two ways to increase any network size: *width* and *depth*. The properties of the width (given a fixed depth) have been extensively studied in the literature: recent work by Yang et al. [25] identified the *Maximal Update Parametrization* (μP) that guarantees maximal feature learning in the infinite width limit.³ Another benefit of μP is hyperparameter transfer which enables hyperparameter tuning on smaller models; the optimal hyperparameter choice for the smaller model remains optimal for larger models (i.e., models with larger width). However, despite the achievements of large-scale deep models and the theoretical understanding of scaling width, increasing the depth of neural networks still has both practical limitations and theoretical difficulties. In practice, increasing depth beyond some level often results in performance degradation and/or significant shifts in the optimal hyperparameters. In theory, unlike increasing width, increasing depth introduces new parameters that significantly change the training dynamics. In this paper, we aim to solve this problem by extending μP to include depth scaling. We call the depth scaling Depth- μP .

*Equal contribution.

†Work partially done at the National University of Singapore.

³Here maximal feature learning refers to $\Theta(1)$ change in features in the infinite width limit. This should be contrasted with the lazy training regime where the change in features is of order $\Theta(n^{-1/2})$.

The issue of depth scaling has persisted over time. A decade ago, deep neural networks experienced significant degradation problems — having more than a few dozen layers would increase the training error instead of improving the model’s performance. This was partly due to the vanishing or exploding gradient problem that affects the efficient propagation of information through the network. The introduction of residual networks (ResNet) [8, 9, 18] has partially resolved this issue, allowing for the training of deeper networks with improved performance. ResNet is constructed by layering *residual blocks*, which are composed of a series of convolutional layers and then an element-wise addition with the input. This element-wise addition (commonly referred to as *skip connection*) is a significant innovation of ResNet and remains an important ingredient in modern architectures including Transformers [19].

Specifically, in a residual architecture, the l -th residual block is formulated as

$$x^l = x^{l-1} + g^l(x^{l-1}; W^l),$$

where x^{l-1} is the input, x^l is the output, W^l are the parameters of the block, and g^l (often called the *residual branch*) is a mapping that defines the layer (e.g. a stack of convolutions in ResNet, or SelfAttention and MLP in a Transformer). In this work, we focus on the case where g^l is a biasless perceptron with (or without) activation.

The stacking of many residual blocks causes an obvious issue even at the initialization — the norm of x^l grows with l , so the last layer features do not have a stable norm when increasing the depth. Intuitively, one can stabilize these features by scaling the residual branches with a depth-dependent constant. However, scaling the residual branches with arbitrarily small constants might result in no feature learning in the large depth limit since the gradients will also be multiplied with the scaling factor. In this paper, we propose Depth- μ P, a principled approach to scaling up the depth of residual networks, enabling the training of arbitrarily deep networks while achieving *feature learning* and maximizing *feature diversity* among nearby layers. Our framework extends the previous results on μ P which deals with optimal width scaling [25]. It completes the width scaling and hence provides a full width and depth scaling recipe that guarantees maximal feature learning and hyperparameter transfer across width and depth. Depth- μ P contains the following modifications to the standard practice:

1. There is a multiplier for each residual branch before adding to its input, which is inversely proportional to the square root of L (where L is the depth). Formally, with a constant a independent from L ,

$$x^l = x^{l-1} + \frac{a}{\sqrt{L}} \cdot g^l(x^{l-1}; W^l). \quad (1)$$

2. We set the learning rate of W^l so that the update of W^l during training is proportional to $1/\sqrt{L}$. We derive different learning rate schemes for different optimization algorithms based on this principle. For Adam, because it is scale-invariant to the gradient, the learning rate of W^l is set to be η/\sqrt{L} . On the other hand, the learning rate of W^l for SGD is set as a constant η because the gradient of W^l is already of size $1/\sqrt{L}$ due to the multiplier.

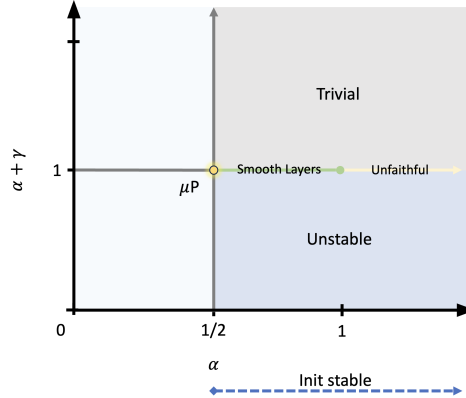


Figure 1: Behaviors of scaling strategies with a branch multiplier $L^{-\alpha}$ and parameter update size proportional to $L^{-\gamma}$.

Figure 1 compares Depth- μ P with general branch multiplier scaling $L^{-\alpha}$ and update size scaling $L^{-\gamma}$ (where Depth- μ P has $\alpha = \gamma = \frac{1}{2}$; see Appendix I for more details), showing the uniqueness of Depth- μ P. Specifically, limiting our analysis to block depth 1 (i.e., g^l is a biasless perceptron, W^l is a single matrix), we show that Depth- μ P leads to the following properties: **a)** At the initialization, each one of the L residual blocks contributes $\Theta(1/\sqrt{L})$ to the main branch. These L contributions are independent of each other, so the sum of them is of size $\Theta(1)$; **b)** During training, the contribution of the update of each residual block is $\Theta(1/L)$ due to the combining effect of the learning rate and multiplier. The contributions of the updates are highly correlated, so they sum up to $\Theta(1)$; **c)** More importantly, we classify all depth limits and show Depth- μ P yields the optimal limit. This implies that the optimal hyperparameters of the networks in Depth- μ P are approximately invariant wrt depth.

With Depth- μ P, we successfully train networks comprising thousands of residual blocks, while also showcasing the transferability of hyperparameters across depth.

While the block depth 1 case admits a positive result, we show that the block depth ≥ 2 case does not and cannot (Appendix M). The basic issue is the weights in different layers within a block is forced to interact additively instead of multiplicatively when depth is large, if one wants to retain diversity. This causes block depth ≥ 2 to have worse performance than block depth 1 and for the optimal hyperparameters to shift with depth. We demonstrate this pedagogically on resnet with MLP blocks but also on Megatron transformer [16] trained on Common Crawl (see Appendix N.3). These observations entail the need to rethink the current approach to hyperparameter transfer.

We refer the reader to Appendix A for a comprehensive literature review.

2 Settings

Notation. Let L be the depth of the network, i.e., the number of residual blocks, and n be the width of the network, i.e. the dimension of all hidden representations x^0, \dots, x^L . Let $\xi \in \mathbb{R}^{d_{\text{in}}}$ be the input of the network, $U \in \mathbb{R}^{n \times d_{\text{in}}}$ be the input layer, and $V \in \mathbb{R}^{n \times e}$ be the output layer, so that $x^0 = U\xi$ and the model output w.r.t. ξ is $f(\xi) \triangleq V^\top x^L$. Let ℓ be the loss function absorbing the label, and δx^l be the gradient of x^l w.r.t. the loss. We denote variables at t -th training step by adding t as a subscript, e.g., the input at step t is ξ_t , the hidden representation of l -th layer at step t is x_t^l , and the model output at step t is f_t . Let T be the number of training steps. We assume the optimizer is Adam and the learning rate is $\eta n^{-c} L^{-\gamma}$.

Setup. We consider an L -hidden-layer residual network with biasless perceptron blocks:

$$\forall l \in [L], \quad h^l = W^l x^{l-1}, \quad x^l = x^{l-1} + aL^{-\alpha} \text{MS}(\phi(h^l)),$$

where $x^0 = U\xi$, the network output $f = V^\top x^L$, and MS refers to Mean Subtraction and is given by $\text{MS}(x) = x - \langle x, \mathbf{1} \rangle / n = Gx$ with $G = I - \mathbf{1}\mathbf{1}^\top / n$. We follow μ P [21] for the *widthwise* scaling, i.e., the initialization of U, V, W^l are i.i.d. zero-mean Gaussian with variance 1, n^{-2}, n^{-1} respectively, and the c in the learning rate $\eta n^{-c} L^{-\gamma}$ for U, V and W^l are 0, 1, 1 respectively, i.e., the learning rate of W^l is $\eta n^{-1} L^{-\gamma}$. In sum, we use (α, γ) pair to decide the *depthwise* scaling of our whole setup. Note MS and μ P are keys to our analysis. In Appendix I, we discuss their necessities.

3 Classification of Depthwise Parametrization

In this section, we establish our main results. We state the results as ‘‘claims’’ instead of theorems. In appendix K.4, we provide ‘‘heuristic’’ proofs that can be made rigorous under non-trivial technical conditions. We believe this additional layer of complexity is unneeded and does not serve the purpose of this paper. For readers who are not familiar with Tensor Programs (TP) framework [23, 25], we provide a warm-up (in Appendix B) for intuition, and rigorously analyze the linear case (in Appendix C) to give a gentle introduction of using TP. We also showcase the correctness of the claims in this section by proving them rigorously in the linear setting in Appendix E.

Before delving into the details, let us first define the (simplified) notions of stability, faithfulness, non-triviality, and feature learning (formal definitions and explanations of claims can be found in Appendix J). Hereafter, all the asymptotic notations such as \mathcal{O}, Ω and o should be understood in the limit ‘‘ $n \rightarrow \infty$, then $L \rightarrow \infty$ ’’.

Stability. We say a parametrization is *stable at step t* if $\mathbf{h}_t^l, \mathbf{x}_t^l = \mathcal{O}(1), \forall l \in [L]$, and $\mathbf{f}_t = \mathcal{O}(1)$.

Faithful. We say a parametrization is *faithful at step t* if $\mathbf{h}_t^l = \Theta(1)$ for all $l \in [L]$.

Nontriviality. We say a parametrization is *trivial* if for any time $t \geq 1$, $\mathbf{f}_t - \mathbf{f}_0 = o(1)$. We say the parametrization is *nontrivial* otherwise.

Feature Learning. We say a parametrization induces *feature learning* if there exist $t \geq 1$, and for any $\lambda > 0$, we have $\Delta \mathbf{h}_t^{\lfloor \lambda L \rfloor} = \Theta(1)$.

Then we have the following claims.

Claim 3.1. *A parametrization is stable at initialization iff $\alpha \geq 1/2$.*

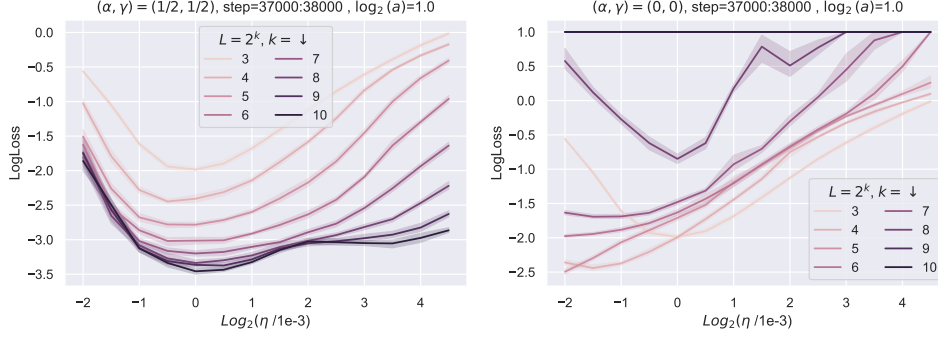


Figure 2: Train logloss versus the learning rate for varying depth with Depth- μ P (Left) and Standard Parametrization without any scaling (Right). We cap the logloss at 1.

Claim 3.2. Consider a parametrization that is stable at initialization. Then the following hold (separately from each other): a) It is stable during training as well iff $\alpha + \gamma \geq 1$; b) It is nontrivial iff $\alpha + \gamma \leq 1$. Therefore, it is both stable and nontrivial iff $\alpha + \gamma = 1$.

Claim 3.3. Consider a stable and nontrivial parametrization. The following hold (separately from each other). a) It is faithful at initialization iff $\alpha \geq 1/2$. As a result, $\alpha = 1/2$ is the minimal choice of α that guarantees faithfulness. b) It is faithful during training iff $\alpha \leq 1$. Therefore, a stable and nontrivial parametrization is faithful iff $\alpha \in [1/2, 1]$.

With the three claims above, we have shrunk the space of “interesting” (α, γ) to $\alpha \in [1/2, 1]$ and $\alpha + \gamma = 1$. To further distinguish these parametrizations, we introduce feature diversity.

Definition 3.1 (Feature Diversity Exponent). We say a parametrization has feature diversity exponent $\kappa \geq 0$ if κ is the maximal value such that for all $\lambda \in [0, 1]$ and sufficiently small $\epsilon > 0$, and all time t ,

$$\frac{1}{\sqrt{n}} \left\| \mathbf{x}_t^{[(\lambda+\epsilon)L]} - \mathbf{x}_t^{[\lambda L]} \right\| = \Omega(\epsilon^{1-\kappa}).$$

We say a parametrization is *redundant* if $\kappa = 0$.

In other words, the feature diversity exponent κ is a measure of how different the outputs are in layers that are close to each other. With $\kappa = 0$, the output of each layer is essentially the same as the output of the previous layer in the sense that the rate of change from one layer to the next is bounded (at least locally), and hence the network is intuitively “wasting” parameters. We provide more discussion on feature diversity in Appendix L, where we predict the better performance of absolute value activation by feature diversity.

Claim 3.4. Consider a stable and nontrivial parametrization that is furthermore faithful during training (but not necessarily at initialization). Then it is redundant ($\kappa = 0$) if $\alpha \in (1/2, 1]$.

Claim 3.5 (Depth- μ P). $\alpha = \gamma = 1/2$ is the unique parametrization that is stable, nontrivial, faithful, induces feature learning, and achieves maximal feature diversity with $\kappa = 1/2$.

4 Experiments

In this section, we provide empirical evidence to show the optimality of Depth- μ P scaling and the transferability of some quantities across depth. We train vanilla residual network with block depth 1 (1 MLP layer in each residual block) on CIFAR-10 dataset using ReLU activation, Adam optimizer, batch size 64, for 50 epochs (input and output layers are fixed). In Figure 2, we show the training loss versus learning rate for depths 2^k , for $k \in \{3, 4, \dots, 10\}$. For Depth- μ P, a convergence pattern can be observed for the optimal learning rate as depth grows. For standard parametrization without any depth scaling ($\alpha = \gamma = 0$), the optimal learning rate exhibits a significant shift as depth grows, and the performance degrades when L grows, suggesting that standard parametrization is not suitable for depth scaling.

In Appendix N, we conduct comprehensive experiments that a) compare Depth- μ P with other parametrization, b) replace MS with LayerNorm, c) verify our linear case analysis, d) verify feature diversity claims, and e) explore Depth- μ P on Transformers.

Acknowledgement

We thank Huishuai Zhang, Jeremy Bernstein, Edward Hu, Michael Santacrose, Liyuan Liu for their helpful comments and discussion. D. Yu was supported by NSF and ONR. Part of this work was done during D. Yu's internship at Microsoft.

References

- [1] Z. Allen-Zhu, Y. Li, and Z. Song. A convergence theory for deep learning via over-parameterization, 2019.
- [2] L. Chizat and F. Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport, 2018.
- [3] L. Chizat, E. Oyallon, and F. Bach. On lazy training in differentiable programming, 2020.
- [4] B. Hanin and D. Rolnick. How to start training: The effect of initialization and architecture, 2018.
- [5] S. Hayou. On the infinite-depth limit of finite-width neural networks, 2023.
- [6] S. Hayou and G. Yang. Width and depth limits commute in residual networks, 2023.
- [7] S. Hayou, E. Clerico, B. He, G. Deligiannidis, A. Doucet, and J. Rousseau. Stable resnet. In A. Banerjee and K. Fukumizu, editors, *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pages 1324–1332. PMLR, 13–15 Apr 2021. URL <https://proceedings.mlr.press/v130/hayou21a.html>.
- [8] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [9] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 630–645. Springer, 2016.
- [10] A. Jacot, F. Gabriel, and C. Hongler. Neural tangent kernel: Convergence and generalization in neural networks, 2020.
- [11] S. Jelassi, B. Hanin, Z. Ji, S. J. Reddi, S. Bhojanapalli, and S. Kumar. Depth dependence of μp learning rates in relu mlps, 2023.
- [12] L. Liu, X. Liu, J. Gao, W. Chen, and J. Han. Understanding the difficulty of training transformers. *arXiv preprint arXiv:2004.08249*, 2020.
- [13] L. Noci, S. Anagnostidis, L. Biggio, A. Orvieto, S. P. Singh, and A. Lucchi. Signal propagation in transformers: Theoretical perspectives and the role of rank collapse, 2022.
- [14] L. Noci, C. Li, M. B. Li, B. He, T. Hofmann, C. Maddison, and D. M. Roy. The shaped transformer: Attention models in the infinite depth-and-width limit, 2023.
- [15] OpenAI. Gpt-4 technical report, 2023.
- [16] M. Shoenybi, M. Patwary, R. Puri, P. LeGresley, J. Casper, and B. Catanzaro. Megatron-lm: Training multi-billion parameter language models using model parallelism. *arXiv preprint arXiv:1909.08053*, 2019.
- [17] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. P. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis. Mastering the game of go with deep neural networks and tree search. *Nature*, 529: 484–489, 2016.
- [18] R. K. Srivastava, K. Greff, and J. Schmidhuber. Highway networks, 2015.
- [19] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need, 2017.
- [20] G. Yang. Scaling limits of wide neural networks with weight sharing: Gaussian process behavior, gradient independence, and neural tangent kernel derivation, 2020.
- [21] G. Yang. Tensor programs ii: Neural tangent kernel for any architecture, 2020.
- [22] G. Yang. Tensor programs i: Wide feedforward or recurrent neural networks of any architecture are gaussian processes, 2021.
- [23] G. Yang and E. J. Hu. Tensor programs iv: Feature learning in infinite-width neural networks. In *International Conference on Machine Learning*, pages 11727–11737. PMLR, 2021.

- [24] G. Yang and E. Littwin. Tensor programs ivb: Adaptive optimization in the infinite-width limit, 2023.
- [25] G. Yang, E. J. Hu, I. Babuschkin, S. Sidor, X. Liu, D. Farhi, N. Ryder, J. Pachocki, W. Chen, and J. Gao. Tensor programs v: Tuning large neural networks via zero-shot hyperparameter transfer. *arXiv preprint arXiv:2203.03466*, 2022.
- [26] H. Zhang, Y. N. Dauphin, and T. Ma. Fixup initialization: Residual learning without normalization, 2019.
- [27] D. Zou, Y. Cao, D. Zhou, and Q. Gu. Stochastic gradient descent optimizes over-parameterized deep relu networks, 2018.

A Related Works

A.1 Width Scaling and μP

The infinite-width limit of neural networks has been a topic of extensive research in the literature. Numerous studies have predominantly focused on examining the behavior of various statistical quantities at initialization. Some works have gone beyond the initialization stage to explore the dynamics of feature learning in neural networks.

Lazy training. With standard parametrization, a learning rate of order $\mathcal{O}(n^{-1})$,⁴ n being the width, yields the so-called lazy training regime in the infinite-width limit, where the features remain roughly constant throughout training [3, 25]. This regime is also known as the Neural Tangent Kernel (NTK) regime and its convergence properties have been extensively studied in the literature [10, 1, 2, 27].

Feature learning and μP . Recent empirical studies (e.g. [25]) have provided compelling evidence that feature learning plays a crucial role in the success of deep learning. It is widely acknowledged that the remarkable performance achieved by deep neural networks can be attributed to their ability to acquire meaningful representations through the process of training. Consequently, scaling the network architecture emerges as a natural choice to enhance the performance of such models.

In this context, μP (Maximal Update Parameterization), introduced in [25], has emerged as a promising approach for maximizing feature learning while simultaneously preventing feature explosion as the network width increases, given a fixed depth. Notably, μP facilitates hyperparameter transfer across varying network widths. This means that instead of tuning hyperparameters directly on large models, one can optimize them on smaller models and utilize the same set of hyperparameters for larger models.

The derivation of μP leverages the Tensor Programs framework [22, 20, 21, 23, 25], which provides valuable tools for capturing the behavior of neural networks in the infinite-width regime during the training process.

A.2 Depth Scaling

While increasing the width of neural networks can lead to improved performance, increasing the depth of the network also yields significant performance gains, and most state-of-the-art models use deep architectures. The introduction of skip connections [8, 9] played a pivotal role in enabling the training of deep networks. However, it became apparent that even with skip connections and normalization layers, training deep networks remains a challenging task [12]. Moreover, tuning hyperparameters for large depth networks is a time-and-resource-consuming task.

To address the challenges associated with training deep networks, several studies have proposed scaling the network blocks using a depth-dependent scaler to ensure stability of features and gradients at initialization [7, 4, 26, 13, 5, 6, 14]. However, these works primarily focus on stability at initialization and lack insights into the dynamics during the training process. For instance, one might argue that features can still experience explosive growth if the learning rate is not properly chosen. Therefore, an effective depth scaling approach should not only ensure stability at initialization but also provide guidelines for scaling the learning rate.

This motivation underlies the development of Depth- μP , which offers a comprehensive framework for depth scaling. Depth- μP encompasses block multipliers and learning rate scaling, providing a complete recipe for training deep networks. In the case of Multi-Layer Perceptrons (MLPs) (no skip connections), Jelassi et al. [11] showed that a learning rate scaling of $depth^{-3/2}$ guarantees stability after the initial gradient step. However, it remains unclear how the learning rate should be adjusted beyond the first step, and this scaling is not suitable for architectures with residual connections.

⁴We also obtain the lazy infinite-width limit with the NTK parametrization and a $\mathcal{O}(n^{-1/2})$ learning rate.

B Warm-Up: An Intuitive Explanation with Linear Networks

Let us begin with a simple example that provides the necessary intuition underpinning our depth scaling strategy. Given a depth L , width n , consider a linear residual network of the form

$$\begin{aligned} x^0 &= U\xi, \\ \forall l \in [L], \quad x^l &= x^{l-1} + \frac{1}{\sqrt{L}}W^l x^{l-1}, \\ f &= V^\top x^L, \end{aligned}$$

where the weight matrices $W^l \in \mathbb{R}^{n \times n}$ and U, V are input and output weight matrices that we assume to be fixed during training.

B.1 Optimal Scaling of the Learning Rate

To simplify the analysis, we consider gradient updates based on a single datapoint. The first gradient step is given by

$$W_1^l = W_0^l - \eta G_0^l,$$

where η is the learning rate, and G_0^l is a matrix with update directions. For instance, we have the following expressions for G_0^l with SGD and Adam:

- SGD: $G_0^l = \frac{1}{\sqrt{L}}\delta x^l \otimes x^{l-1}$, where $\delta x^l \stackrel{\text{def}}{=} \frac{\partial \ell}{\partial x^l}$ for some loss function ℓ .⁵
- Adam⁶: $G_0^l = \text{sign}\left(\frac{1}{\sqrt{L}}\delta x^l \otimes x^{l-1}\right)$.

In both cases, δx^l and x^{l-1} are computed for a single training datapoint ξ_0 . The last layer features x^L (for some input ξ) are given by $x^L = \prod_{l=1}^L \left(I + \frac{1}{\sqrt{L}}W^l\right) x^0$.⁷ We use the subscript t to refer to training step. After the first gradient step, we have the following

$$x_1^L = \prod_{l=1}^L \left(I + \frac{1}{\sqrt{L}}W_1^l\right) x^0 = x_0^L - \frac{\eta}{\sqrt{L}}A_L + \mathcal{O}(\eta^2), \quad (2)$$

where $A_L = \sum_{l=1}^L \left[\prod_{k>l} \left(I + \frac{1}{\sqrt{L}}W_0^k\right)\right] G_0^l \left[\prod_{k<l} \left(I + \frac{1}{\sqrt{L}}W_0^k\right)\right] x^0$. We argue that A_L behaves as $\Theta(L)$ (in L_2 norm). This is due to the $1/\sqrt{L}$ scaling factor. To see this, we further simplify the analysis by considering the case $d_{in} = n = d_{out} = 1$ (single neuron per layer) and the squared loss. In this case, the term A_L simplifies to

$$A_L = \sum_{l=1}^L \prod_{k \neq l} \left(1 + \frac{1}{\sqrt{L}}W_0^k\right) G_0^l x_0.$$

Scaling for SGD. With SGD, we have that $G_0^l = \frac{1}{\sqrt{L}} \prod_{k \neq l} \left(1 + \frac{1}{\sqrt{L}}W_0^k\right) x_0 \delta x^L$, where $\delta x^L = (Vx^L - y(\xi_0))$ and $y(\xi_0)$ is the target output. Therefore, it is easy to see that

$$\mathbb{E}A_L^2 = \frac{1}{L} \mathbb{E} \left(\sum_{l=1}^L \prod_{k \neq l} \left(1 + \frac{1}{\sqrt{L}}W_0^k\right) \delta x^L x_0^2 \right)^2 = \Theta\left(\frac{1}{L}L^2\right) = \Theta(L),$$

where we have used the fact that $\mathbb{E} \left(1 + \frac{1}{\sqrt{L}}W_0^k\right)^{2p} = 1 + \Theta(L^{-1})$, for any positive integer p .

⁵We use δ for gradient because we want to distinguish from d in depth differential equations that appear later in the paper.

⁶For the sake of simplification, we consider SignSGD in this section, which can be seen as a memory-less version of Adam. The analysis is valid for any training algorithm that gives $\Theta(1)$ gradients.

⁷To avoid any confusion, here we define the matrix product by $\prod_{l=1}^L A_l = A_L \times A_{L-1} \cdots \times A_1$.

Hence, the magnitude of the first order term in eq. (2) is given by

$$\mathbb{E} \left[\left(\frac{\eta}{\sqrt{L}} A_l \right)^2 \right] = \Theta(\eta^2),$$

which shows that the update is stable in depth as long as $\eta = \Theta(1)$ in depth. More precisely, this is the maximal choice of learning rate that does not lead to exploding features as depth increases.

Scaling for Adam. With Adam, we have $G_0^l = \pm 1$, and therefore we obtain

$$\mathbb{E} A_l^2 = \mathbb{E} \left(\sum_{l=1}^L \prod_{k \neq l} \left(1 + \frac{1}{\sqrt{L}} W_0^k \right) x_0 \right)^2 = \Theta(L^2),$$

where we have used the same arguments as before. In this case, the first order term in eq. (2) is given by

$$\mathbb{E} \left[\left(\frac{\eta}{\sqrt{L}} A_l \right)^2 \right] = \Theta(\eta^2 L^{-1}).$$

Therefore, the maximal learning rate that one can choose without exploding the features is given by $\eta = \Theta(L^{-1/2})$.

Summary: By ensuring that parameter update is $\Theta(1/\sqrt{L})$, the features remain stable while feature update is $\Theta(1)$. This $\Theta(1)$ update is due to the accumulation of $\Theta(1/L)$ correlated terms across depth.

B.2 Convergence when Depth goes to ∞

Let us look at x_1^L again in the simple case $d_{in} = d_{out} = n = 1$ and analyze its behaviour when $L \rightarrow \infty$. This paragraph is only intended to give an intuition for the convergence. A rigorous proof of such convergence will be later presented in the paper. Let us consider the case with SGD training with learning rate $\eta = 1$ and let $M_{L,l} = \prod_{k \neq l} \left(1 + \frac{1}{\sqrt{L}} W_0^k \right)$ and $\tau = (Vx_0^L - y(\xi_0))x^0$. With this, we have the following

$$x_1^L = \prod_{l=1}^L \left(1 + \frac{1}{\sqrt{L}} W_0^l - \frac{1}{L} \tau M_{L,l} \right) x^0. \quad (3)$$

WLOG, let us assume that $x_0^0 > 0$. Then, with high probability (the event that $W_0^l \ll \sqrt{L}$, for some notion of “ \ll ”, occurs with a probability of at least $1 - e^{-L^\alpha}$ for some $\alpha > 0$)⁸, we have that $x_1^L > 0$. We can therefore look at $\log(x_1^L)$ which simplifies the task. Taking the log and using Taylor expansion under a high probability event, we obtain

$$\begin{aligned} \log(x_1^L/x^0) &= \frac{1}{\sqrt{L}} \sum_{l=1}^L W_0^l - \frac{1}{L} \sum_{l=1}^L \tau M_{L,l} + \frac{\sum_{l=1}^L (W_0^l)^2}{L} + \mathcal{O}(L^{-1+\epsilon}) \\ &= \frac{1}{\sqrt{L}} \sum_{l=1}^L W_0^l - \tau x_0^L \frac{1}{L} \sum_{l=1}^L \frac{1}{1 + \frac{1}{\sqrt{L}} W_0^l} + \frac{\sum_{l=1}^L (W_0^l)^2}{L} + \mathcal{O}(L^{-1+\epsilon}), \end{aligned}$$

for some $\epsilon > 0$. The first and third terms $\frac{1}{\sqrt{L}} \sum_{l=1}^L W_0^l$ and $\frac{\sum_{l=1}^L (W_0^l)^2}{L}$ converge (almost surely) to a standard Gaussian and 1, respectively. The second term also converges naturally, since x_0^L converges in L_2 to a Log-Normal random variable ([5]) and with a delicate treatment (involving high probability bounds), one can show that the term $\frac{1}{L} \sum_{l=1}^L \frac{1}{1 + \frac{1}{\sqrt{L}} W_0^l}$ converges (in L_2 norm) at large depth. This

implies that one should expect x_1^L to have some notion of weak convergence as depth grows. Note that the same analysis becomes much more complicated for general width $n > 0$. To avoid dealing with high probability bounds, a convenient method consists of taking the width to infinity first $n \rightarrow \infty$, then analyzing what happens as depth increases. We discuss this in the next section.

⁸This follows from simple concentration inequalities for sub-exponential random variables.

B.3 A Discussion on the General Case

Difficulty of generalizing to the nonlinear case. The extension to the general width scenario ($n > 1$) necessitates a more intricate treatment of the term A_l to find optimal scaling rules, yet the proposed scaling remains optimal for general width. This preliminary analysis lays the groundwork for proposing a specific learning rate scaling scheme that maximizes feature learning. Moreover, demonstrating the optimality of this scaling strategy in the presence of non-linearities is a non-trivial task. The primary challenge stems from the correlation among the post-activations induced during the training process. Overcoming these challenges requires a rigorous framework capable of addressing the large depth limit of crucial quantities in the network.

For this purpose, we employ the Tensor Program framework to investigate the behavior of essential network quantities in the infinite-width-then-depth limit. By leveraging this framework, our theoretical findings establish that the aforementioned scaling strategy remains optimal for general networks with skip connections. Our framework considers the setup where the width is taken to infinity first, followed by depth. This represents the case where $1 \ll \text{depth} \ll \text{width}$, which encompasses most practical settings (e.g. Large Language Models).

The critical role of Initialization. A naive approach to depth scaling can be as follows: since the weights W_t^k might become highly correlated during training, one has to scale the blocks with $1/L$. To understand this, let us assume a block multiplier of $L^{-\alpha}$ and consider the scenario of perfect correlation where all weights are equal, i.e., $W_t^k = W$ for every $k \in 1, \dots, L$. In this case, the last layer features can be expressed as $x^L = (I + L^{-\alpha}W)^L x_0$. When $\alpha = 1/2$, the features are likely to exhibit an explosive growth with increasing depth, while opting for $\alpha = 1$ is guaranteed to stabilize the features.

However, in this paper, we demonstrate that this intuition does not align with practical observations. Contrary to expectations, the features do not undergo an explosive growth as the depth increases when $\alpha = 1/2$. This phenomenon is attributed to two key factors: random initialization and learning rate scaling with depth. These factors ensure that the weight matrices never become highly correlated in this particular fashion during the training process.

In summary, while a naive depth scaling strategy based on scaling blocks might suggest the need for $\alpha = 1$ to stabilize the features, our findings reveal that in practice, this is not the case. The interplay of random initialization and learning rate scaling effectively prevents the features from experiencing explosive growth, even with the choice of $\alpha = 1/2$.

C SGD Training Dynamics of Infinitely Deep Linear Networks

In this section, we continue to study the linear neural network with residual connections under Depth- μ P. Using the Tensor Program framework [24], we rigorously derive the training dynamics of SGD for the linear residual network when the width and the depth sequentially go to infinity. The road map of our analysis consists the following three steps.

1. We first take the width of the network to infinity by the Tensor Program framework [24]. As a result, instead of tracking vectors and matrices along the training trajectory, we track random variables that correspond to the vectors, that is, for a vector $x \in \mathbb{R}^n$ that appears in the computation of the training, the coordinates of x can be viewed as iid copies of random variable $\llbracket x \rrbracket$ (called a *ket*) when $n \rightarrow \infty$.⁹
2. Since the network is linear, every random variable can be written as a linear combination of a set of zero-mean “base” random variables by the Master Theorem of Tensor Programs [24]. Therefore, we can track the random variables by analyzing the coefficients of their corresponding linear combinations, along with the covariance between the “base” random variables.
3. Since the number of random variables and the number of “base” random variables scale linearly with L , the coefficients of all random variables can be represented by a six-

⁹The definition of $\llbracket x \rrbracket$ requires the coordinates of x is $\mathcal{O}(1)$ w.r.t. n , and $\llbracket x \rrbracket$ is trivial if the coordinates of x is $o(1)$ w.r.t. n . Therefore, for x whose coordinates are not $\Theta(1)$, we normalize x by multiplying polynomial of n so the resulting vector has coordinates $\Theta(1)$.

dimensional tensor, where two of the dimensions have shape L . We then map the tensor to a set of functions whose input domain is $[0, 1] \times [0, 1]$. Finally, we claim that the functions converge when $L \rightarrow \infty$, and identify their limits as the solution of a set of functional integrals.

In Appendix N.1, we conduct a thorough empirical verification of our theory in the linear case. The experiments clearly show the convergence of deep linear residual networks under Depth- μ P.

Assumptions and Notations Recall the linear network is given by

$$\begin{aligned} x^0 &= U\xi, \\ \forall l \in [L], \quad x^l &= \frac{a}{\sqrt{L}} W^l x^{l-1} + x^{l-1}, \\ f &= V^\top x^L. \end{aligned}$$

For convenience, we assume $a = 1$, the SGD learning rate of W^l is 1. We add t as a subscript to any notation to denote the same object but at t -th training step, e.g., the input at step t is a single datapoint ξ_t , the hidden output of l -th layer at step t is x_t^l , and the model output at step t is f_t . Let T be the number of training steps. Let ℓ_t be the loss function absorbing the label at time t , and χ_t be the derivative of the loss at time t , i.e., $\chi_t = \ell'_t(f_t)$. Let $\delta x_t^l = \partial \ell_t / \partial x_t^l$, and $\tilde{\delta} x_t^l = n \delta x_t^l$ is the normalized version of δx_t^l .

The Tensor Program analysis heavily depends on the scaling of initialization and learning rate of U, V, W w.r.t n . In this paper, we use μ P as the scaling w.r.t. n since it maximizes feature learning in the large width limit [23]. Without loss of generality, we follow [23] and assume the input and output dimension is 1, i.e., $\xi \in \mathbb{R}, f \in \mathbb{R}$. For a clean presentation, we additionally assume U, V are frozen during training in this section and each coordinate of W is initialized with i.i.d. Gaussian of variance $1/n$.

C.1 Width Limit under μ P

As the first step, we take the width of the network n to infinity using Tensor Programs (TP). As briefly mentioned in the road map of the section, the TP framework characterizes each vector involved in the training procedure by a random variable when $n \rightarrow \infty$. For a vector $x \in \mathbb{R}^n$ that has roughly iid coordinates, we write $\|x\rangle \in \mathbb{R}$ (called a *ket*) to denote a random variable such that x 's entries look like iid copies of $\|x\rangle$. Then for any two vector $x, y \in \mathbb{R}^n$ that have roughly iid coordinates, their limiting inner product by n can be written as $\lim_{n \rightarrow \infty} \frac{x^\top y}{n} = \mathbb{E} \|x\rangle \cdot \|y\rangle$, which we write succinctly as $\langle x \| y \rangle$. Deep linear network with SGD is a simple example for this conversion from vectors to random variables. As shown in Program 1, we define a series of scalars (f_t° and χ_t°) and random variables ($\|U\rangle, \|nV\rangle, \|x_t^l\rangle, \|\delta x_t^l\rangle, \|W_t^l x_t^{l-1}\rangle, \|W_t^{l\top} \delta x_t^l\rangle$) using the ket notations. For better understanding, we provide a brief introduction to TP below.

Tensor Programs (TP) in a nutshell. When training a neural network, one can think of this procedure as a process of successively creating new vectors and scalars from an initial set of random vectors and matrices (initialization weights), and some deterministic quantities (dataset in this case). In the first step, the forward propagation creates the features x_0^l where the subscript 0 refers to initialization, and the scalar f_0 , which is the network output. In the first backward pass, the output derivative χ_0 is computed, then the gradients δx_0^l are backpropagated. (Since the coordinates of δx_0^l vanish to 0 when $n \rightarrow \infty$, TP instead tracks its normalized version $\tilde{\delta} x_0^l \stackrel{\text{def}}{=} n \cdot \delta x_0^l$.) New vectors are created and appended to the TP as training progresses. When the width n goes to infinity, vectors of size n in the TP (e.g., the features x_t^l , and normalized gradients $\tilde{\delta} x_t^l$) see their coordinates converge to roughly iid random variables (e.g., $\|x_t^l\rangle$ and $\|\tilde{\delta} x_t^l\rangle$ in Program 1), and other scalar quantities (e.g., f_t and χ_t) converge to deterministic values (e.g., f_t° and χ_t° in Program 1) under proper parametrization (μ P). The Master Theorem [25] captures the behaviour of these quantities by characterizing the *infinite-width* limit of the training process. For more in-depth definitions and details about TP, we refer the reader to [25].

Program 1: Random Variables induced from Tensor Program for the Linear Network with LR $\eta = 1$ and frozen U, V .

Initial random variables: $\llbracket U \rrbracket, \llbracket nV \rrbracket$ are independent standard Gaussian.

for $t = 0, \dots, T - 1$ **do**

$\llbracket x_t^0 \rrbracket \stackrel{\text{def}}{=} \xi_t \llbracket U \rrbracket;$

for $l = 1, \dots, L$ **do**

$\llbracket W_t^l x_t^{l-1} \rrbracket \stackrel{\text{def}}{=} \llbracket W_0^l x_t^{l-1} \rrbracket - \frac{1}{\sqrt{L}} \sum_{s=0}^{t-1} \llbracket \tilde{\delta} x_s^l \rrbracket \langle x_s^{l-1} \llbracket x_t^{l-1} \rrbracket \rangle;$

$\llbracket x_t^l \rrbracket \stackrel{\text{def}}{=} \llbracket x_t^{l-1} \rrbracket + \frac{1}{\sqrt{L}} \llbracket W_t^l x_t^{l-1} \rrbracket;$

end

$\mathring{f}_t \stackrel{\text{def}}{=} \langle x_t^L \llbracket nV \rrbracket \rangle;$

$\mathring{\chi}_t \stackrel{\text{def}}{=} \ell_t(\mathring{f}_t);$

$\llbracket \delta x_t^L \rrbracket \stackrel{\text{def}}{=} \mathring{\chi}_t \llbracket nV \rrbracket;$

for $l = L, \dots, 1$ **do**

$\llbracket W_t^{l\top} \tilde{\delta} x_t^l \rrbracket \stackrel{\text{def}}{=} \llbracket W_0^{l\top} \tilde{\delta} x_t^l \rrbracket - \frac{1}{\sqrt{L}} \sum_{s=0}^{t-1} \llbracket x_s^{l-1} \rrbracket \langle \tilde{\delta} x_s^l \llbracket \tilde{\delta} x_t^l \rrbracket \rangle;$

$\llbracket \tilde{\delta} x_t^{l-1} \rrbracket \stackrel{\text{def}}{=} \llbracket \tilde{\delta} x_t^l \rrbracket + \frac{1}{\sqrt{L}} \llbracket W_t^{l\top} \tilde{\delta} x_t^l \rrbracket;$

end

end

where $\llbracket W_0^l x_t^{l-1} \rrbracket$ and $\llbracket W_0^{l\top} \tilde{\delta} x_t^l \rrbracket$ are defined in Definition C.1.

Now when we look back to Program 1, the definitions of scalars and random variables should be clear (except for $\llbracket W_0^l x_t^{l-1} \rrbracket$ and $\llbracket W_0^{l\top} \tilde{\delta} x_t^l \rrbracket$). One can find straightforward correspondence between those and their finite counterpart, for example:

- \mathring{f}_t corresponds to f_t , and $\mathring{\chi}_t$ corresponds to χ_t ;
- $\llbracket x_t^l \rrbracket$ corresponds to x_t^l and $\llbracket \tilde{\delta} x_t^l \rrbracket$ corresponds to $\tilde{\delta} x_t^l$. (Recall $\tilde{\delta} x_t^l = n \cdot \delta x_t^l$ is the normalized version of δx_t^l .)
- By SGD, $W_t^l = W_0^l - \frac{1}{\sqrt{L}} \sum_{s=0}^{t-1} \delta x_s^l \otimes x_s^{l-1}$, which corresponds to $\llbracket W_t^l x_t^{l-1} \rrbracket = \llbracket W_0^l x_t^{l-1} \rrbracket - \frac{1}{\sqrt{L}} \sum_{s=0}^{t-1} \llbracket \tilde{\delta} x_s^l \rrbracket \langle x_s^{l-1} \llbracket x_t^{l-1} \rrbracket \rangle$.

Now we can dive into the definition of $\llbracket W_0^l x_t^{l-1} \rrbracket$ and $\llbracket W_0^{l\top} \tilde{\delta} x_t^l \rrbracket$. Let \mathcal{W} be the set of initial random matrices of size $n \times n$, i.e., $\{W_0^1, \dots, W_0^L\}$, and $\mathcal{W}^\top \stackrel{\text{def}}{=} \{W^\top : W \in \mathcal{W}\}$. Let \mathcal{V}_W denote the set of all vectors in training of the form Wy for some y . Then for every $W \in \mathcal{W} \cup \mathcal{W}^\top$, and $Wy \in \mathcal{V}_W$, we can decompose $\llbracket Wy \rrbracket$ into the sum of $\llbracket W\hat{y} \rrbracket$ and $\llbracket Wy \rrbracket$, where $\llbracket W\hat{y} \rrbracket$ is a random variable that act as if W were independent of y , and $\llbracket Wy \rrbracket$ is the random variable capturing the correlation part between W and y . Specifically, let us briefly track what happens to $W_0^l x_t^{l-1}$ during training. In the first step, we have $W_0^l x_0^{l-1}$ which has roughly Gaussian coordinates (in the large width limit). In this case, we have $\llbracket W_0^l x_0^{l-1} \rrbracket = 0$. After the first backprop, we have $\delta x_0^{l-1} = \delta x_0^l + \frac{1}{\sqrt{L}} W_0^{l\top} \delta x_0^l$, which means that the update in W^{l-1} will contain a term of the form $W_0^{l\top} z$ for some vector z . This implies that $W_0^l x_1^{l-1}$ will contain a term of the form $W_0^l W_0^{l\top} z'$ for some vector z' . This term induces an additional correlation term that appears when we take the width to infinity. The $\llbracket W_0^l x_1^{l-1} \rrbracket$ is defined by isolating this additional correlation term from $W_0^l W_0^{l\top} z'$. The remaining term is Gaussian in the infinite-width limit, which defines the term $\llbracket W_0^l x_1^{l-1} \hat{\rrbracket}$. Formally, we present the following definition.

Definition C.1. We define $\llbracket Wy \rrbracket \stackrel{\text{def}}{=} \llbracket W\hat{y} \rrbracket + \llbracket Wy \rrbracket$ for every $W \in \mathcal{W} \cup \mathcal{W}^\top$ and $Wy \in \mathcal{V}_W$, where

- $\llbracket W\hat{y} \rrbracket$ is a Gaussian variable with zero mean. $\forall W \in \mathcal{W} \cup \mathcal{W}^\top, Wy, Wz \in \mathcal{V}_W$,

$$\text{Cov}(\llbracket W\hat{y} \rrbracket, \llbracket W\hat{z} \rrbracket) \stackrel{\text{def}}{=} \langle y \rrbracket z \rangle.$$

$\forall W, W' \in \mathcal{W} \cup \mathcal{W}^\top, Wy \in \mathcal{V}_W, W'z \in \mathcal{V}_{W'}, \llbracket Wy \rrbracket$ and $\llbracket W'z \rrbracket$ are independent if $W \neq W'$. $\llbracket Wy \rrbracket$ is also independent from $\llbracket U \rrbracket$ and $\llbracket nV \rrbracket$.

- $\llbracket Wy \rrbracket$ is defined to be a linear combination of $\{\llbracket z \rrbracket : W^\top z \in \mathcal{V}_{W^\top}\}$. Then we can unwind any $\llbracket y \rrbracket$ inductively as a linear combination of $\llbracket \bullet \rrbracket, \llbracket U \rrbracket$ and $\llbracket nV \rrbracket$, which allows us to fully define

$$\llbracket Wy \rrbracket \stackrel{\text{def}}{=} \sum_{W^\top z \in \mathcal{V}_{W^\top}} \llbracket z \rrbracket \cdot \frac{\partial \llbracket y \rrbracket}{\partial \llbracket W^\top z \rrbracket}.$$

C.2 Depthwise Scaling of Random Variables

As mentioned in Definition C.1, both $\llbracket x_t^l \rrbracket$ and $\llbracket \tilde{\delta} x_t^{l-1} \rrbracket$ can be written as linear combination of “base” random variables: $\{\llbracket W_0^m x_s^{m-1} \rrbracket\}_{s \in \{0, \dots, t\}, m \in [L]}, \{\llbracket W_0^{m^\top} \tilde{\delta} x_s^m \rrbracket\}_{s \in \{0, \dots, t\}, m \in [L]}, \llbracket U \rrbracket$ and $\llbracket nV \rrbracket$. Moreover, the coefficients of the linear combinations can be calculated in a recursive way: by expanding $\llbracket W_0^l x_t^{l-1} \rrbracket$ using Definition C.1, we have

$$\llbracket x_t^l \rrbracket = \llbracket x_t^{l-1} \rrbracket + \frac{1}{\sqrt{L}} \llbracket W_0^l x_t^{l-1} \rrbracket + \frac{1}{\sqrt{L}} \sum_{s=1}^{t-1} \llbracket \tilde{\delta} x_s^l \rrbracket \left(\frac{\partial \llbracket x_t^{l-1} \rrbracket}{\partial \llbracket W_0^{l^\top} \tilde{\delta} x_s^l \rrbracket} - \frac{1}{\sqrt{L}} \langle x_s^{l-1} \llbracket x_t^{l-1} \rrbracket \rangle \right).$$

The recursive formula for $\llbracket \tilde{\delta} x_t^l \rrbracket$ is similar.

Using this induction, we claim in the linear combinations, the coefficient of every $\llbracket \bullet \rrbracket$ is $\mathcal{O}(1/\sqrt{L})$, and the coefficient of $\llbracket U \rrbracket$ and $\llbracket nV \rrbracket$ is $\mathcal{O}(1)$. We also claim the covariance between any pairs of random variables in the form of $\llbracket x_t^l \rrbracket$ and $\llbracket \tilde{\delta} x_t^{l-1} \rrbracket$ is $\mathcal{O}(1)$.

Proposition C.2. $\forall t, \forall s \leq t, \forall l, m, \forall \llbracket y \rrbracket \in \{\llbracket x_t^l \rrbracket, \llbracket \tilde{\delta} x_t^l \rrbracket\}$,

$$\frac{\partial \llbracket y \rrbracket}{\partial \llbracket W_0^m x_s^{m-1} \rrbracket} = \mathcal{O}\left(\frac{1}{\sqrt{L}}\right), \frac{\partial \llbracket y \rrbracket}{\partial \llbracket W_0^{m^\top} \tilde{\delta} x_s^m \rrbracket} = \mathcal{O}\left(\frac{1}{\sqrt{L}}\right), \frac{\partial \llbracket y \rrbracket}{\partial \llbracket U \rrbracket} = \mathcal{O}(1), \frac{\partial \llbracket y \rrbracket}{\partial \llbracket nV \rrbracket} = \mathcal{O}(1).$$

$\forall t, s, l, m, \forall \llbracket y \rrbracket \in \{\llbracket x_t^l \rrbracket, \llbracket \tilde{\delta} x_t^l \rrbracket\}, \forall \llbracket z \rrbracket \in \{\llbracket x_s^m \rrbracket, \llbracket \tilde{\delta} x_s^m \rrbracket\}$,

$$\langle y \llbracket z \rrbracket \rangle = \mathcal{O}(1).$$

The reasoning of Proposition C.2 is provided in Appendix D. Note the computation of covariance can also be written as a recursive formula. The reasoning relies essentially on an inductive argument.

C.3 Infinite Depth Limit

Now we formalize our argument above and obtain the formula describing the dynamics of the network when $L \rightarrow \infty$. We first write the coefficients of the linear combinations as a six dimensional tensor $\Gamma_{t,s,a,b,l,m}$, where $t, s \in \{0, \dots, T-1\}, a, b \in \{0, 1\}, l, m \in [L]$. Specifically, $\Gamma_{t,s,a,b,l,m}$ represents the derivative of $\llbracket x_t^l \rrbracket$ and $\llbracket \tilde{\delta} x_t^l \rrbracket$ w.r.t. $\llbracket W_0^m x_s^{m-1} \rrbracket$ and $\llbracket W_0^{m^\top} \tilde{\delta} x_s^m \rrbracket$. Here, we use 0 to denote kets appears in the forward pass ($\llbracket x_t^l \rrbracket$ and $\llbracket W_0^m x_s^{m-1} \rrbracket$), and 1 to denote kets in the backward pass ($\llbracket \tilde{\delta} x_t^l \rrbracket$ and $\llbracket W_0^{m^\top} \tilde{\delta} x_s^m \rrbracket$). Formally, $\Gamma_{t,s,0,0,l,m} = \frac{\partial \llbracket x_t^l \rrbracket}{\partial \llbracket W_0^m x_s^{m-1} \rrbracket}, \Gamma_{t,s,0,1,l,m} = \frac{\partial \llbracket x_t^l \rrbracket}{\partial \llbracket W_0^{m^\top} \tilde{\delta} x_s^m \rrbracket},$

$$\Gamma_{t,s,1,0,l,m} = \frac{\partial \llbracket \tilde{\delta} x_t^l \rrbracket}{\partial \llbracket W_0^m x_s^{m-1} \rrbracket}, \Gamma_{t,s,1,1,l,m} = \frac{\partial \llbracket \tilde{\delta} x_t^l \rrbracket}{\partial \llbracket W_0^{m^\top} \tilde{\delta} x_s^m \rrbracket}.$$

However, it is hard to describe the limit of Γ because its size increases along with L . Therefore, we define the following set of functions $\{\Gamma_{t,s,a,b} : [0, 1] \times (0, 1] \rightarrow \mathbb{R}\}_{t \in \{0, \dots, T-1\}, s \in \{-1, \dots, t\}, a, b \in \{0, 1\}}$. For $s \geq 0$,

$$\Gamma_{t,s,a,b}(p, q) = \sqrt{L} \cdot \Gamma_{t,s,a,b,[Lp],[Lq]}$$

$$\text{For } s = -1, \Gamma_{t,-1,0,0}(p, q) = \frac{\partial \llbracket x_t^{[Lp]} \rrbracket}{\partial \llbracket U \rrbracket}, \Gamma_{t,-1,0,1}(p, q) = \frac{\partial \llbracket x_t^{[Lp]} \rrbracket}{\partial \llbracket nV \rrbracket}, \Gamma_{t,-1,1,0}(p, q) = \frac{\partial \llbracket \tilde{\delta} x_t^{[Lp]} \rrbracket}{\partial \llbracket U \rrbracket}, \Gamma_{t,-1,1,1}(p, q) = \frac{\partial \llbracket \tilde{\delta} x_t^{[Lp]} \rrbracket}{\partial \llbracket nV \rrbracket}.$$

Here l, m are normalized to $[0, 1]$ so the input domain of Γ 's are identical for different L ; $\Gamma_{t,s,a,b,l,m}$ is multiplied by \sqrt{L} because $\Gamma_{t,s,a,b,l,m} = \mathcal{O}(1/\sqrt{L})$ by Proposition C.2; and the extra $s = -1$ case helps us also capture the derivative w.r.t. $\|U\rangle$ and $\|nV\rangle$.

Similarly, we can also define another set of function $\{C_{t,s,a} : (0, 1] \rightarrow \mathbb{R}\}_{t,s \in \{-1, \dots, T-1\}, a \in \{0,1\}}$ to describe the covariance between the ‘‘base’’ random variables: $\forall p \in (0, 1]$, let $l = \lceil Lp \rceil$,

$$\begin{aligned} \bullet C_{t,s,0}(p) &\stackrel{\text{def}}{=} \text{Cov}(\|W_0^l x_t^{l-1}\rangle, \|W_0^l x_s^{l-1}\rangle) = \langle x_t^{l-1} \| x_s^{l-1} \rangle, \\ \bullet C_{t,s,1}(p) &\stackrel{\text{def}}{=} \text{Cov}(\|W_0^{l\top} \tilde{\delta} x_t^l\rangle, \|W_0^{l\top} \tilde{\delta} x_s^l\rangle) = \langle \tilde{\delta} x_t^l \| \tilde{\delta} x_s^l \rangle, \end{aligned}$$

For $t = -1$, $C_{-1,-1,0}(p) \stackrel{\text{def}}{=} \text{Cov}(\|U\rangle, \|U\rangle) = 1$, and $C_{-1,-1,1}(p) \stackrel{\text{def}}{=} \text{Cov}(\|nV\rangle, \|nV\rangle) = 1$. By Definition C.1, the ‘‘base’’ random variables of different ‘‘groups’’ are independent, so we only tracks the covariance listed above.

Using this definition of Γ and C , it is convenient to write their recursive formula in the following lemma.

Lemma C.3 (Finite depth recursive formula for Γ and C (Informal version of Lemma D.1)). *Γ and C can be computed recursively as follows:*

$$\begin{aligned} \Gamma_{t,r,0,b} \left(\frac{l}{L}, q \right) &= \Gamma_{t,r,0,b} \left(\frac{l-1}{L}, q \right) + \mathbb{I}_{[(t=r) \wedge (b=0) \wedge (l=\lceil Lq \rceil)]} \\ &\quad + \frac{1}{L} \sum_{s=0}^{t-1} \Gamma_{s,r,1,b} \left(\frac{l}{L}, q \right) \left(\Gamma_{t,s,0,1} \left(\frac{l-1}{L}, \frac{l}{L} \right) - C_{t,s,0} \left(\frac{l}{L} \right) \right). \end{aligned}$$

$$\begin{aligned} \Gamma_{t,r,1,b} \left(\frac{l-1}{L}, q \right) &= \Gamma_{t,r,1,b} \left(\frac{l}{L}, q \right) + \mathbb{I}_{[(t=r) \wedge (b=1) \wedge (l=\lceil Lq \rceil)]} \\ &\quad + \frac{1}{L} \sum_{s=0}^{t-1} \Gamma_{s,r,0,b} \left(\frac{l-1}{L}, q \right) \left(\Gamma_{t,s,1,0} \left(\frac{l}{L}, \frac{l}{L} \right) - C_{t,s,1} \left(\frac{l}{L} \right) \right). \end{aligned}$$

$$C_{t,s,a}(p) = \sum_{t'=-1}^t \sum_{s'=-1}^s \sum_{b \in \{0,1\}} \int_0^1 \Gamma_{t,t',a,b}(l/L, q) C_{t',s',b}(q) \Gamma_{s,s',a,b}(l/L, q) dq,$$

where $l = \lceil Lp \rceil - 1$ if $a = 0$, and $l = \lceil Lp \rceil$ if $a = 1$.

The proof of Lemma C.3 is straightforward from Program 1. In Appendix D, we also give a formal proof that Γ and C converge when L grows to infinity, in the case where L is powers of 2. The restriction on L being powers of 2 is imposed for the convenience of the proof, and the convergence of Γ and C is true in the general case. Moreover, we derive the infinite depth behavior based on the recursion of Γ and C in Lemma C.3.

Proposition C.4 (Infinite depth limit of Γ and C (Informal version of Proposition D.2)). *In the limit $L \rightarrow \infty$, we have*

$$\begin{aligned} \Gamma_{t,r,0,b}(p, q) &= \mathbb{I}_{[(t=r) \wedge (b=0) \wedge (p \geq q)]} + \int_0^p \sum_{s=0}^{t-1} \Gamma_{s,r,1,b}(p', q) \cdot (\Gamma_{t,s,0,1}(p', p') - C_{t,s,0}(p')) dp'; \\ \Gamma_{t,r,1,b}(p, q) &= \mathbb{I}_{[(t=r) \wedge (b=1) \wedge (p \leq q)]} + \int_p^1 \sum_{s=0}^{t-1} \Gamma_{s,r,0,b}(p', q) \cdot (\Gamma_{t,s,1,0}(p', p') - C_{t,s,1}(p')) dp'; \\ C_{t,s,a}(p) &= \sum_{t'=-1}^t \sum_{s'=-1}^s \sum_{b \in \{0,1\}} \int_0^1 \Gamma_{t,t',a,b}(p, q) C_{t',s',b}(q) \Gamma_{s,s',a,b}(p, q) dq. \end{aligned}$$

The proof of Proposition C.4 follows from Lemma C.3. A rigorous proof requires first showing the existence of a solution of the integral functional satisfied by the couple (Γ, C) . The solution is typically a fixed point of the integral functional in Proposition C.4. After showing the existence, one needs to show that (Γ, C) converges to this limit. This typically requires controlling the difference

between finite-depth and infinite-depth solutions and involves obtaining upper-bounds on error propagation. The existence is guaranteed under mild conditions on the integral functional. We omit here the full proof for existence and assume that the functional is sufficiently well-behaved for this convergence result to hold. The formal proof of the convergence of Γ and C for $L = 2^k$ ($k \in \mathbb{N}$) in Appendix D is a showcase of the correctness of the proposition.

D Details of the linear case

D.1 Proof sketch of Proposition C.2

Here we provide a proof sketch of Proposition C.2, the formal prove is implied by the existence of Γ and C in the infinite depth limit.

Proof sketch. The claims can be reasoned by induction on t and l . Let us take $\llbracket x_t^l \rrbracket$ as an example, since $\llbracket \tilde{\delta} x_t^{l-1} \rrbracket$ is symmetric with $\llbracket x_t^l \rrbracket$. By expanding the definition of $\llbracket x_t^l \rrbracket$, we have

$$\llbracket x_t^l \rrbracket = \llbracket x_t^{l-1} \rrbracket + \frac{1}{\sqrt{L}} \llbracket W_0^l x_t^{l-1} \widehat{\cdot} \rrbracket + \frac{1}{\sqrt{L}} \sum_{s=1}^{t-1} \llbracket \tilde{\delta} x_s^l \rrbracket \left(\frac{\partial \llbracket x_t^{l-1} \rrbracket}{\partial \llbracket W_0^{l\top} \tilde{\delta} x_s^l \rrbracket} - \frac{1}{\sqrt{L}} \langle x_s^{l-1} \llbracket x_t^{l-1} \rrbracket \rangle \right).$$

Note by induction, $\langle x_s^{l-1} \llbracket x_t^{l-1} \rrbracket \rangle = \mathcal{O}(1)$ and $\frac{\partial x_t^{l-1}}{\partial \llbracket W_0^{l\top} \tilde{\delta} x_s^l \rrbracket} = \mathcal{O}(1/\sqrt{L})$, so

$$\begin{aligned} \llbracket x_t^l \rrbracket &= \llbracket x_t^{l-1} \rrbracket + \frac{1}{\sqrt{L}} \llbracket W_0^l x_t^{l-1} \widehat{\cdot} \rrbracket + \mathcal{O} \left(\frac{1}{L} \right) \sum_{s=1}^{t-1} \llbracket \tilde{\delta} x_s^l \rrbracket \\ &= \xi_t \llbracket U \rrbracket + \sum_{m=1}^l \frac{1}{\sqrt{L}} \llbracket W_0^m x_t^{m-1} \widehat{\cdot} \rrbracket + \mathcal{O} \left(\frac{1}{L} \right) \sum_{m'=1}^l \sum_{s'=1}^{t-1} \llbracket \tilde{\delta} x_{s'}^{m'} \rrbracket. \end{aligned}$$

Then by unwinding $\llbracket \tilde{\delta} x_{s'}^{m'} \rrbracket$ and noting that by induction, $\forall s < t$, $\frac{\partial \llbracket \tilde{\delta} x_{s'}^{m'} \rrbracket}{\partial \llbracket W_0^{m'} x_s^{m'-1} \rrbracket} = \mathcal{O} \left(\frac{1}{\sqrt{L}} \right)$, $\frac{\partial \llbracket \tilde{\delta} x_{s'}^{m'} \rrbracket}{\partial \llbracket W_0^{m\top} \tilde{\delta} x_s^{m'} \rrbracket} = \mathcal{O} \left(\frac{1}{\sqrt{L}} \right)$, $\frac{\partial \llbracket \tilde{\delta} x_{s'}^{m'} \rrbracket}{\partial \llbracket U \rrbracket} = \mathcal{O}(1)$, $\frac{\partial \llbracket \tilde{\delta} x_{s'}^{m'} \rrbracket}{\partial \llbracket nV \rrbracket} = \mathcal{O}(1)$, we have

$$\frac{\partial \llbracket x_t^l \rrbracket}{\partial \llbracket W_0^m x_s^{m-1} \rrbracket} = \mathcal{O} \left(\frac{1}{\sqrt{L}} \right), \frac{\partial \llbracket x_t^l \rrbracket}{\partial \llbracket W_0^{m\top} \tilde{\delta} x_s^{m'} \rrbracket} = \mathcal{O} \left(\frac{1}{\sqrt{L}} \right), \frac{\partial \llbracket x_t^l \rrbracket}{\partial \llbracket U \rrbracket} = \mathcal{O}(1), \frac{\partial \llbracket x_t^l \rrbracket}{\partial \llbracket nV \rrbracket} = \mathcal{O}(1).$$

Also by unwinding, $\forall \llbracket y \rrbracket \in \{ \llbracket x_s^m \rrbracket, \llbracket \tilde{\delta} x_s^m \rrbracket \}$,

$$\begin{aligned} \langle y \llbracket x_t^l \rrbracket \rangle &= \sum_{m'} \sum_{s'} \sum_{t'} \frac{\partial \llbracket x_t^l \rrbracket}{\partial \llbracket W_0^{m'} x_{t'}^{m'-1} \rrbracket} \cdot \frac{\partial \llbracket y \rrbracket}{\partial \llbracket W_0^{m'} x_{t'}^{m'-1} \rrbracket} \cdot \langle x_{t'}^{m'-1} \llbracket x_{s'}^{m'-1} \rrbracket \rangle \\ &\quad + \sum_{m'} \sum_{s'} \sum_{t'} \frac{\partial \llbracket x_t^l \rrbracket}{\partial \llbracket W_0^{m'\top} \tilde{\delta} x_{t'}^{m'} \rrbracket} \cdot \frac{\partial \llbracket y \rrbracket}{\partial \llbracket W_0^{m'\top} \tilde{\delta} x_{t'}^{m'} \rrbracket} \cdot \langle \tilde{\delta} x_{t'}^{m'} \llbracket \tilde{\delta} x_{s'}^{m'} \rrbracket \rangle \\ &\quad + \frac{\partial \llbracket x_t^l \rrbracket}{\partial \llbracket U \rrbracket} \cdot \frac{\partial \llbracket y \rrbracket}{\partial \llbracket U \rrbracket} + \frac{\partial \llbracket x_t^l \rrbracket}{\partial \llbracket nV \rrbracket} \cdot \frac{\partial \llbracket y \rrbracket}{\partial \llbracket nV \rrbracket} \\ &= \mathcal{O}(1). \end{aligned} \quad \square$$

D.2 Formal recursive formula of Γ and C

By the same way of expanding $\llbracket x_t^l \rrbracket$ and $\langle y \llbracket x_t^l \rrbracket \rangle$, we formally derive the recursive formula for Γ and C below.

Lemma D.1 (Finite depth recursive formula for Γ and C). Γ can be computed recursively as follows:

For $t = 0, \dots, T-1$,

$$\bullet \forall q \in (0, 1], \Gamma_{t,-1,0,q}(0, q) = \xi_t,$$

- For $l = 1, \dots, L$, $\forall r \leq t$, $\forall p \in (\frac{l-1}{L}, \frac{l}{L}]$, $\forall q \in (0, 1]$, $\forall b \in \{0, 1\}$,

$$C_{t,s,0}(p) = \sum_{t'=-1}^t \sum_{s'=-1}^s \sum_{b \in \{0,1\}} \int_0^1 \Gamma_{t,t',0,b} \left(\frac{l-1}{L}, q \right) C_{t',s',b}(q) \Gamma_{s,s',0,b} \left(\frac{l-1}{L}, q \right) dq;$$

$$\begin{aligned} \Gamma_{t,r,0,b}(p, q) &= \Gamma_{t,r,0,b} \left(\frac{l-1}{L}, q \right) + \mathbb{I}_{[(t=r) \wedge (b=0) \wedge (l=\lceil Lq \rceil)]} \\ &\quad + \frac{1}{L} \sum_{s=0}^{t-1} \Gamma_{s,r,1,b} \left(\frac{l}{L}, q \right) \left(\Gamma_{t,s,0,1} \left(\frac{l-1}{L}, \frac{l}{L} \right) - C_{t,s,0} \left(\frac{l}{L} \right) \right). \end{aligned}$$

- $\mathring{f}_t = \Gamma_{t,-1,0,1}(1, 1)$,
- $\mathring{\chi}_t = \ell'_t(\mathring{f}_t)$,
- $\forall q \in (0, 1], \Gamma_{t,-1,1,1}(1, q) = \mathring{\chi}_t$,
- For $l = L, \dots, 1$, $\forall r \leq t$, $\forall p \in (\frac{l-2}{L}, \frac{l-1}{L}]$, $\forall q \in (0, 1]$, $\forall b \in \{0, 1\}$,

$$C_{t,s,1} \left(p + \frac{1}{L} \right) = \sum_{t'=-1}^t \sum_{s'=-1}^s \sum_{b \in \{0,1\}} \int_0^1 \Gamma_{t,t',1,b}(l/L, q) C_{t',s',b}(q) \Gamma_{s,s',1,b}(l/L, q) dq;$$

$$\begin{aligned} \Gamma_{t,r,1,b}(p, q) &= \Gamma_{t,r,1,b} \left(\frac{l}{L}, q \right) + \mathbb{I}_{[(t=r) \wedge (b=1) \wedge (l=\lceil Lq \rceil)]} \\ &\quad + \frac{1}{L} \sum_{s=0}^{t-1} \Gamma_{s,r,0,b} \left(\frac{l-1}{L}, q \right) \left(\Gamma_{t,s,1,0} \left(\frac{l}{L}, \frac{l}{L} \right) - C_{t,s,1} \left(\frac{l}{L} \right) \right). \end{aligned}$$

The proof is straightforward from Program 1. The recursive nature of Γ and C yields the following infinite-depth behavior.

Proposition D.2 (Infinite depth limit of Γ and C). *In the limit $L \rightarrow \infty$, we have $\forall p \in [0, 1], q \in (0, 1], b \in \{0, 1\}$:*

$$\Gamma_{t,-1,0,0}(0, q) = \xi_t;$$

$$\Gamma_{t,r,0,b}(p, q) = \mathbb{I}_{[(t=r) \wedge (b=0) \wedge (p \geq q)]} + \int_0^p \sum_{s=0}^{t-1} \Gamma_{s,r,1,b}(p', q) \cdot (\Gamma_{t,s,0,1}(p', p') - C_{t,s,0}(p')) dp';$$

$$\mathring{f}_t = \Gamma_{t,-1,0,1}(1, 1);$$

$$\mathring{\chi}_t = \ell'_t(\mathring{f}_t);$$

$$\Gamma_{t,-1,1,1}(1, q) = \mathring{\chi}_t;$$

$$\Gamma_{t,r,1,b}(p, q) = \mathbb{I}_{[(t=r) \wedge (b=1) \wedge (p \leq q)]} + \int_p^1 \sum_{s=0}^{t-1} \Gamma_{s,r,0,b}(p', q) \cdot (\Gamma_{t,s,1,0}(p', p') - C_{t,s,1}(p')) dp';$$

$$C_{t,s,a}(p) = \sum_{t'=-1}^t \sum_{s'=-1}^s \sum_{b \in \{0,1\}} \int_0^1 \Gamma_{t,t',a,b}(p, q) C_{t',s',b}(q) \Gamma_{s,s',a,b}(p, q) dq.$$

D.3 Convergence of Γ and C when $L = 2^k$

In this section, we prove Γ and C will converge when $L \rightarrow \infty$. For convenience, we will only consider the case when $L = 2^k$ for some integer k . To distinguish Γ and C corresponding to different L , we add the depth as the superscript, i.e., Γ^L and C^L .

Theorem D.3. $\forall t \leq T, s < t, a \in \{0, 1\}, b \in \{0, 1\}, \forall p \in [0, 1], q \in (0, 1]$,

- $\{\Gamma_{t,s,a,b}^{2^k}(p, q)\}_{k \in \mathbb{N}}$ is a Cauchy sequence,

- $\{C_{t,s,a}^{2^k}(p)\}_{k \in \mathbb{N}}$ is a Cauchy sequence.

The proof is by induction on t . We will prove the following claims (A) (B) (C) (D) on $t > 0$ given they are satisfied for any $s < t$. For $t = 0$, (A) (B) (C) (D) are trivial.

Assumption on $s < t$ Assume $\exists c > 1$ such that $\forall L > L'$ and $L = 2^k$ for $k \in \mathbb{N}$, $\forall s < t$, $\forall r < s$,

$$(A) \quad \forall p \in \{0, \frac{1}{L}, \dots, 1\}, q \in (0, 1],$$

$$|\Gamma_{s,r,a,b}^{L/2}(p, q) - \Gamma_{s,r,a,b}^L(p, q)| \leq c/L, \quad |C_{s,r,a}^{L/2}(p, q) - C_{s,r,a}^L(p, q)| \leq c/L.$$

$$(B) \quad |\Gamma_{s,r,a,b}^L(p, q)| \leq c, |C_{s,r,a}^L(p)| \leq c.$$

$$(C) \quad C_{s,r,a}^L(p) \text{ is } c\text{-Lipschitz w.r.t. } p, \text{ and } \Gamma_{s,r,a,b}^L(p, q) \text{ is } c\text{-Lipschitz w.r.t. } p.$$

$$(D) \quad |\Gamma_{s,r,0,1}^L(p - \frac{1}{L}, p + \frac{1}{L}) - \Gamma_{s,r,0,1}^L(p - \frac{1}{L}, p)| \leq c/L, |\Gamma_{s,r,1,0}^L(p, p) - \Gamma_{s,r,1,0}^L(p, p - \frac{1}{L})| \leq c/L.$$

Remark (A) indicates that $\{\Gamma_{s,r,a,b}^{2^k}\}_k$ and $\{C_{s,r,a}^{2^k}\}_k$ converge. We only care about $r < s$ because $C_{s,s,a}^L$ will never be used, and $\Gamma_{s,s,a,b}^L$ is known: for $p \in \{0, \frac{1}{L}, \dots, 1\}$,

$$\Gamma_{s,s,a,b}^L(p, q) = \mathbb{I}[(a = 0) \wedge (b = 0) \wedge (p \geq q)] + \mathbb{I}[(a = 1) \wedge (b = 1) \wedge (p + 1/L \leq q)].$$

Proof for t -th step (the forward pass) In the following subsections, we will prove inductively on increasing order of all $L > L'$ and $L = 2^k$, and increasing order of $p \in \{0, 1/L, \dots, 1\}$ that $\forall s < t$,

$$(D0) \quad |\Gamma_{t,s,0,1}^L(p, p + \frac{2}{L}) - \Gamma_{t,s,0,1}^L(p, p + \frac{1}{L})| \leq c_2 \exp(c_1 p)/L;$$

$$(C0) \quad \text{For } s < t, |\Gamma_{t,s,0,b}^L(p, q) - \Gamma_{t,s,0,b}^L(p - \frac{1}{L}, q)| \leq t c c_2 \exp(c_1(p - \frac{1}{L}))/L;$$

$$(B0) \quad |\Gamma_{t,s,0,b}^L(p, q)| \leq c_2 \exp(c_1(p - \frac{1}{2L}));$$

$$(A0) \quad |\Gamma_{t,s,0,b}^{L/2}(p, q) - \Gamma_{t,s,0,b}^L(p, q)| \leq c_3 c_2 \exp(c_1(p - \frac{1}{2L}))/L;$$

$$(C1) \quad |C_{t,s,0}^L(p + \frac{1}{L}) - C_{t,s,0}^L(p)| \leq c_4 c_2 \exp(c_1(p - \frac{1}{L}))/L;$$

$$(B1) \quad |C_{t,s,0}^L(p + \frac{1}{L})| \leq c_2 \exp(c_1 p);$$

$$(A1) \quad |C_{t,s,0}^{L/2}(p + \frac{1}{L}) - C_{t,s,0}^L(p + \frac{1}{L})| \leq c_5 c_2 \exp(c_1 p)/L,$$

where $c_2 = \max\{\xi_t^2, |\xi_t|\} \exp(c_1/2L')$, $c_3 = 3ct$, $c_4 = 4t(t+1)c^2 + 2tc$, $c_5 = c_4 + 1$, $c_1 = c^3 t(4ct + 2c_4 + 29) + tc(3c_4 + 14) + c(2c_4 + 2c)$.

Proof for t -th step (the backward pass) Similar bounds also apply to $\Gamma_{t,s,1,b}$ and $C_{t,s,1}$ by induction on decreasing order of p .

Conclusion Combining both backward pass and forward pass at time t shows (A)(B)(C)(D) also hold for $s = t$ with a larger (but constant) c . Thus, (A)(B)(C)(D) hold for any constant s by induction on training steps.

D.3.1 $\Gamma_{t,s,0,b}^L(p, q)$ in forward pass (Proof for D0, C0, B0, A0)

We first consider

$$\begin{aligned} \Gamma_{t,r,0,b}^L(p, q) &= \Gamma_{t,r,0,b}^L\left(p - \frac{1}{L}, q\right) + \mathbb{I}[(t = r) \wedge (b = 0) \wedge (Lp = \lceil Lq \rceil)] \\ &\quad + \frac{1}{L} \sum_{s=0}^{t-1} \Gamma_{s,r,1,b}^L(p, q) \left(\Gamma_{t,s,0,1}^L\left(p - \frac{1}{L}, p\right) - C_{t,s,0}^L(p) \right). \end{aligned}$$

(D0) Difference between $\Gamma_{t,s,0,1}^L(p, p + \frac{2}{L})$ and $\Gamma_{t,s,0,1}^L(p, p + \frac{1}{L})$ Assume $p \geq 1/L$ ($p = 0$ is trivial), let $q = p+1/L, q' = p+2/L$, note that $\Gamma_{s,r,1,b}^L(p, q) = \Gamma_{s,r,1,b}^L(p, q')$ since $p+1/L \leq q \leq q'$, so for $r < t$,

$$\begin{aligned}
& |\Gamma_{t,r,0,b}^L(p, q) - \Gamma_{t,r,0,b}^L(p, q')| \\
& \leq |\Gamma_{t,r,0,b}^L(p - \frac{1}{L}, q) - \Gamma_{t,r,0,b}^L(p - \frac{1}{L}, q')| \\
& \quad + \frac{1}{L} \sum_{s=0}^{t-1} |\Gamma_{s,r,1,b}^L(p, q) - \Gamma_{s,r,1,b}^L(p, q')| \left| \Gamma_{t,s,0,1}^L(p - \frac{1}{L}, p) - C_{t,s,0}^L(p) \right| \\
& \leq c_2 \exp(c_1(p - \frac{1}{L}))/L + \frac{1}{L} \cdot t \cdot c/L \cdot 2c_2 \exp(c_1(p - \frac{1}{L})) \\
& = (1 + 2ct/L)c_2 \exp(c_1(p - \frac{1}{L}))/L \leq c_2 \exp(c_1 p)/L,
\end{aligned}$$

as $c_1 \geq 2ct$.

(C0) Lipschitz w.r.t. p For $r < t$,

$$\begin{aligned}
& |\Gamma_{t,r,0,b}^L(p, q) - \Gamma_{t,r,0,b}^L(p - \frac{1}{L}, q)| \\
& = \left| \frac{1}{L} \sum_{s=0}^{t-1} \Gamma_{s,r,1,b}^L(p, q) \left(\Gamma_{t,s,0,1}^L(p - \frac{1}{L}, p) - C_{t,s,0}^L(p) \right) \right| \\
& \leq \frac{1}{L} \sum_{s=0}^{t-1} c \left(c_2 \exp(c_1(p - \frac{1}{L})) + c_2 \exp(c_1(p - \frac{1}{L})) \right) \\
& = ct c_2 \exp(c_1(p - \frac{1}{L}))/L.
\end{aligned}$$

(B0) Bounded Again assume $p \geq 1/L$ ($p = 0$ is trivial because $c_2 \geq |\xi_t| \exp(c_1/2L)$), since $|\Gamma_{t,r,0,b}^L(p - \frac{1}{L}, q)| \leq c_2 \exp(c_1(p - \frac{1}{L}))$, we can bound $|\Gamma_{t,r,0,b}^L(p, q)|$:

$$\begin{aligned}
|\Gamma_{t,r,0,b}^L(p, q)| & \leq c_2 \exp(c_1(p - \frac{1}{L})) + ct c_2 \exp(c_1(p - \frac{1}{L}))/L \\
& = c_2 \exp(c_1(p - \frac{1}{L}))(1 + ct/L) \\
& \leq c_2 \exp(c_1(p - \frac{1}{2L})),
\end{aligned}$$

as long as $c_1 \geq 2ct$.

(A0) Difference between L and $L/2$ bounded When $p = 0$, it is trivial. When $p = 1/L$, it is also trivial by Lipschitz w.r.t. p , which results

$$|\Gamma_{t,r,0,b}^{L/2}(p, q) - \Gamma_{t,r,0,b}^L(p, q)| \leq 3ctc_2/L \leq c_3c_2 \exp(c_1/2L)/L.$$

When $p \geq 2/L$, since

$$\Gamma_{t,r,0,b}^{L/2}(p, q) = \Gamma_{t,r,0,b}^{L/2}(p - \frac{2}{L}, q) + \frac{2}{L} \sum_{s=0}^{t-1} \Gamma_{s,r,1,b}^{L/2}(p, q) \left(\Gamma_{t,s,0,1}^{L/2}(p - \frac{2}{L}, p) - C_{t,s,0}^{L/2}(p) \right),$$

we compare it with $\Gamma_{t,r,0,b}^L(p, q)$ expanded based on its previous two steps

$$\begin{aligned}
\Gamma_{t,r,0,b}^L(p, q) & = \Gamma_{t,r,0,b}^L(p - \frac{2}{L}, q) + \frac{1}{L} \sum_{s=0}^{t-1} \Gamma_{s,r,1,b}^L(p, q) \left(\Gamma_{t,s,0,1}^L(p - \frac{1}{L}, p) - C_{t,s,0}^L(p) \right) \\
& \quad + \frac{1}{L} \sum_{s=0}^{t-1} \Gamma_{s,r,1,b}^L(p - \frac{1}{L}, q) \left(\Gamma_{t,s,0,1}^L(p - \frac{2}{L}, p - \frac{1}{L}) - C_{t,s,0}^L(p - \frac{1}{L}) \right).
\end{aligned}$$

In order to bridge the two above, namely matching the inputs for Γ and C , we need a middle term

$$\tilde{\Gamma}_{t,r,0,b}^L(p, q) = \Gamma_{t,r,0,b}^L\left(p - \frac{2}{L}, q\right) + \frac{2}{L} \sum_{s=0}^{t-1} \Gamma_{s,r,1,b}^L(p, q) \left(\Gamma_{t,s,0,1}^L\left(p - \frac{2}{L}, p\right) - C_{t,s,0}^L(p) \right).$$

Now we can bound $|\Gamma_{t,r,0,b}^L(p, q) - \tilde{\Gamma}_{t,r,0,b}^L(p, q)|$, and $|\tilde{\Gamma}_{t,r,0,b}^L(p, q) - \Gamma_{t,r,0,b}^{L/2}(p, q)|$ separately, which add up to be the bound for $|\Gamma_{t,r,0,b}^L(p, q) - \Gamma_{t,r,0,b}^{L/2}(p, q)|$.

$$\begin{aligned} & |\Gamma_{t,r,0,b}^L(p, q) - \tilde{\Gamma}_{t,r,0,b}^L(p, q)| \\ & \leq \frac{1}{L} \sum_{s=0}^{t-1} |\Gamma_{s,r,1,b}^L(p, q)| \left| \Gamma_{t,s,0,1}^L\left(p - \frac{1}{L}, p\right) - \Gamma_{t,s,0,1}^L\left(p - \frac{2}{L}, p\right) \right| \\ & \quad + \frac{1}{L} \sum_{s=0}^{t-1} \left| \Gamma_{s,r,1,b}^L\left(p - \frac{1}{L}, q\right) \left(\Gamma_{t,s,0,1}^L\left(p - \frac{2}{L}, p - \frac{1}{L}\right) - C_{t,s,0}^L\left(p - \frac{1}{L}\right) \right) \right. \\ & \quad \quad \left. - \Gamma_{s,r,1,b}^L(p, q) \left(\Gamma_{t,s,0,1}^L\left(p - \frac{2}{L}, p\right) - C_{t,s,0}^L(p) \right) \right| \\ & \leq \frac{1}{L} \cdot ct \cdot ctc_2 \exp(c_1(p - \frac{2}{L}))/L + \frac{1}{L} \cdot 2t \cdot c/L \cdot c_2 \exp(c_1(p - \frac{2}{L})) \\ & \quad + \frac{1}{L} \cdot ct \cdot c_2 \exp(c_1(p - \frac{2}{L}))/L + \frac{1}{L} \cdot ct \cdot c_4 c_2 \exp(c_1(p - \frac{2}{L}))/L \\ & = \frac{c^2 t^2 + 3ct + c_4 ct}{L^2} \cdot c_2 \exp(c_1(p - \frac{2}{L})). \end{aligned}$$

and

$$\begin{aligned} & |\Gamma_{t,r,0,b}^{L/2}(p, q) - \tilde{\Gamma}_{t,r,0,b}^L(p, q)| \\ & \leq |\Gamma_{t,r,0,b}^{L/2}\left(p - \frac{2}{L}, q\right) - \Gamma_{t,r,0,b}^L\left(p - \frac{2}{L}, q\right)| \\ & \quad + \frac{1}{L} \sum_{s=0}^{t-1} c \left| \Gamma_{t,s,0,1}^{L/2}\left(p - \frac{2}{L}, p\right) - \Gamma_{t,s,0,1}^L\left(p - \frac{2}{L}, p\right) - C_{t,s,0}^{L/2}(p) + C_{t,s,0}^L(p) \right| \\ & \quad + \frac{1}{L} \sum_{s=0}^{t-1} \frac{c}{L} \left(\left| \Gamma_{t,s,0,1}^L\left(p - \frac{2}{L}, p\right) \right| + \left| C_{t,s,0}^L(p) \right| \right) \\ & \leq \frac{1}{L} \cdot (c_3 c_2 \exp(c_1(p - \frac{1}{L})) + ct(c_3 + c_5) c_2 \exp(c_1(p - \frac{1}{L}))/L + 2t \cdot \frac{c}{L} \cdot c_2 \exp(c_1(p - \frac{1}{L}))) \\ & \leq \frac{c_3 + ct(c_3 + c_5 + 2)/L}{L} \cdot c_2 \exp(c_1(p - \frac{1}{L})). \end{aligned}$$

In sum, as $c_1 \geq \frac{2(c_3 + c_5 + ct + c_4 + 5)}{3}$,

$$\begin{aligned} |\Gamma_{t,r,0,b}^{L/2}(p, q) - \Gamma_{t,r,0,b}^L(p, q)| & \leq \frac{c_3 + ct(c_3 + c_5 + ct + c_4 + 5)/L}{L} c_2 \exp(c_1(p - \frac{1}{L})) \\ & \leq c_2 \exp(c_1(p - \frac{1}{2L}))/L. \end{aligned}$$

D.3.2 $C_{t,s,0}(p + \frac{1}{L})$ in forward pass (Proof for C1, B1, A1)

Now consider $C_{t,s,0}^L$. By expanding

$$C_{t,s,0}^L(p + \frac{1}{L}) = \sum_{t'=-1}^t \sum_{s'=-1}^s \sum_{b \in \{0,1\}} \int_0^1 \Gamma_{t,t',0,b}^L(p, q) C_{t',s',b}^L(q) \Gamma_{s,s',0,b}^L(p, q) dq,$$

we will have

$$\begin{aligned} C_{t,s,0}^L(p + \frac{1}{L}) &= \sum_{t'=-1}^{t-1} \sum_{s'=-1}^s \sum_{b \in \{0,1\}} \int_0^1 \Gamma_{t,t',0,b}^L(p, q) C_{t',s',b}^L(q) \Gamma_{s,s',0,b}^L(p, q) dq \\ &\quad + \sum_{s'=0}^s \int_0^p C_{t,s',0}^L(q) \Gamma_{s,s',0,0}^L(p, q) dq. \end{aligned}$$

(C1) Lipschitz Since $C_{t',s',b}^L$ and $\Gamma_{s,s',0,b}^L$ are bounded and Lipschitz,

$$\begin{aligned} &|C_{t,s,0}^L(p + \frac{1}{L}) - C_{t,s,0}^L(p)| \\ &\leq \sum_{t'=-1}^{t-1} \sum_{s'=-1}^s \sum_{b \in \{0,1\}} \int_0^1 |\Gamma_{t,t',0,b}^L(p, q) - \Gamma_{t,t',0,b}^L(p - \frac{1}{L}, q)| \cdot c^2 dq \\ &\quad + \sum_{t'=-1}^{t-1} \sum_{s'=-1}^s \sum_{b \in \{0,1\}} \int_0^1 |\Gamma_{t,t',0,b}^L(p - \frac{1}{L}, q)| \cdot c \cdot \frac{c}{L} dq \\ &\quad + \sum_{s'=0}^s \frac{1}{L} \cdot |C_{t,s',0}^L(p) \Gamma_{s,s',0,0}^L(p, p)| \\ &\quad + \sum_{s'=0}^s \int_0^{p-\frac{1}{L}} |C_{t,s',0}^L(q)| \cdot \frac{c}{L} dq. \\ &\leq 1/L \cdot (2t(s+1) \cdot ct c_2 \exp(c_1(p - \frac{1}{L})) \cdot c^2 + 2t(s+1) c_2 \exp(c_1(p - \frac{1}{L})) \cdot c^2 \\ &\quad + s \cdot c_2 \exp(c_1(p - \frac{1}{L})) \cdot c + s \cdot c_2 \exp(c_1(p - \frac{1}{L})) \cdot c) \\ &= (4t(s+1)c^2 + 2sc)/L \cdot c_2 \exp(c_1(p - \frac{1}{L})) \\ &\leq c_4 c_2 \exp(c_1(p - \frac{1}{L}))/L. \end{aligned}$$

(B1) Bounded Since $|C_{t,s,0}^L(p)| \leq c_2 \exp(c_1(p - \frac{1}{L}))$, we bound $C_{t,s,0}^L(p + \frac{1}{L})$ as:

$$|C_{t,s,0}^L(p + \frac{1}{L})| \leq c_2 \exp(c_1(p - \frac{1}{L})) \cdot (1 + c_4/L) \leq c_2 \exp(c_1 p),$$

as long as $c_1 \geq c_4$.

(A1) Difference between L and $L/2$ bounded It is easy to see that for $p = 0$,

$$C_{t,s,0}^L(p + \frac{1}{L}) - C_{t,s,0}^{L/2}(p + \frac{1}{L}) = 0,$$

we will prove that for $p \in \{2/L, 4/L, \dots, 1\}$,

$$|C_{t,s,0}^L(p + \frac{1}{L}) - C_{t,s,0}^{L/2}(p + \frac{1}{L})| \leq c_2 \exp(c_1 p)/L.$$

Then by (C1), for $p \in \{1/L, 3/L, \dots, 1 - 1/L\}$,

$$\begin{aligned} |C_{t,s,0}^L(p + \frac{1}{L}) - C_{t,s,0}^{L/2}(p + \frac{1}{L})| &\leq c_2 \exp(c_1(p - \frac{1}{L}))/L + c_4 c_2 \exp(c_1(p - \frac{1}{L}))/L \\ &\leq (c_4 + 1) c_2 \exp(c_1 p)/L \\ &= c_5 c_2 \exp(c_1 p)/L. \end{aligned}$$

Suppose $p \in \{2/L, 4/L, \dots, 1\}$, we compare $C_{t,s,0}^L(p + \frac{1}{L}) - C_{t,s,0}^L(p - \frac{1}{L})$ and $C_{t,s,0}^{L/2}(p + \frac{1}{L}) - C_{t,s,0}^{L/2}(p - \frac{1}{L})$. Intuitively, both of them are $\mathcal{O}(1/L)$, and their difference is $\mathcal{O}(1/L^2)$. In particular,

both of them can be written into four parts:

$$C_{t,s,0}^L(p + \frac{1}{L}) - C_{t,s,0}^L(p - \frac{1}{L}) = \sum_{t'=-1}^{t-1} \sum_{s'=-1}^s \sum_{b \in \{0,1\}} \int_0^1 \left(\Gamma_{t,t',0,b}^L(p, q) - \Gamma_{t,t',0,b}^L(p - \frac{2}{L}, q) \right) C_{t',s',b}^L(q) \Gamma_{s,s',0,b}^L(p, q) dq \quad (\mathcal{E}_1^L)$$

$$+ \sum_{t'=-1}^{t-1} \sum_{s'=-1}^s \sum_{b \in \{0,1\}} \int_0^1 \Gamma_{t,t',0,b}^L(p - \frac{2}{L}, q) C_{t',s',b}^L(q) \left(\Gamma_{s,s',0,b}^L(p, q) - \Gamma_{s,s',0,b}^L(p - \frac{2}{L}, q) \right) dq \quad (\mathcal{E}_2^L)$$

$$+ \sum_{s'=0}^s \int_{p-\frac{2}{L}}^p C_{t,s',0}^L(q) \Gamma_{s,s',0,0}^L(p, q) dq \quad (\mathcal{E}_3^L)$$

$$+ \sum_{s'=0}^s \int_0^{p-\frac{2}{L}} C_{t,s',0}^L(q) (\Gamma_{s,s',0,0}^L(p, q) - \Gamma_{s,s',0,0}^L(p - \frac{2}{L}, q)) dq \quad (\mathcal{E}_4^L)$$

and $C_{t,s,0}^{L/2}(p + \frac{1}{L}) - C_{t,s,0}^{L/2}(p - \frac{1}{L}) = \mathcal{E}_1^{L/2} + \mathcal{E}_2^{L/2} + \mathcal{E}_3^{L/2} + \mathcal{E}_4^{L/2}$ where $\mathcal{E}_i^{L/2}$ is defined in the same way as \mathcal{E}_i^L but with $C^{L/2}$ and $\Gamma^{L/2}$ instead of C^L and Γ^L . Next we bound $|\mathcal{E}_i^L - \mathcal{E}_i^{L/2}|$ one by one:

1. The only hard part to bound in $|\mathcal{E}_i^L - \mathcal{E}_i^{L/2}|$ is

$$|\Gamma_{t,t',0,b}^L(p, q) - \Gamma_{t,t',0,b}^L(p - \frac{2}{L}, q) - (\Gamma_{t,t',0,b}^{L/2}(p, q) - \Gamma_{t,t',0,b}^{L/2}(p - \frac{2}{L}, q))|.$$

By almost the same proof of (A0),

$$\begin{aligned} & |\Gamma_{t,t',0,b}^L(p, q) - \Gamma_{t,t',0,b}^L(p - \frac{2}{L}, q) - (\Gamma_{t,t',0,b}^{L/2}(p, q) - \Gamma_{t,t',0,b}^{L/2}(p - \frac{2}{L}, q))| \\ & \leq \frac{ct(c_3 + c_5 + ct + c_4 + 5)}{L^2} c_2 \exp(c_1(p - \frac{1}{L})). \end{aligned}$$

Then we have

$$\begin{aligned} & |\mathcal{E}_1^L - \mathcal{E}_1^{L/2}| / (2t(s+1)) \\ & \leq \frac{ct(c_3 + c_5 + ct + c_4 + 5)}{L^2} c_2 \exp(c_1(p - \frac{1}{L})) \cdot c \cdot c \\ & \quad + 4ctc_2 \exp(c_1(p - \frac{1}{L})) / L \cdot c / L \cdot c \\ & \quad + 4ctc_2 \exp(c_1(p - \frac{1}{L})) / L \cdot c \cdot c / L \\ & \leq \frac{c^3 t(c_3 + c_5 + ct + c_4 + 13)}{L^2} c_2 \exp(c_1(p - \frac{1}{L})) \end{aligned}$$

2. Bounding $|\mathcal{E}_2^L - \mathcal{E}_2^{L/2}|$ is similar to $|\mathcal{E}_1^L - \mathcal{E}_1^{L/2}|$, where we first bound

$$|\Gamma_{s,s',0,b}^L(p, q) - \Gamma_{s,s',0,b}^L(p - \frac{2}{L}, q) - (\Gamma_{s,s',0,b}^{L/2}(p, q) - \Gamma_{s,s',0,b}^{L/2}(p - \frac{2}{L}, q))| \leq 9c^2 t / L^2.$$

Then we have

$$\begin{aligned} & |\mathcal{E}_2^L - \mathcal{E}_2^{L/2}| / (2t(s+1)) \\ & \leq c_3 c_2 \exp(c_1(p - \frac{2}{L})) / L \cdot c \cdot 2c / L \\ & \quad + c_2 \exp(c_1(p - \frac{2}{L})) \cdot c / L \cdot 2c / L \\ & \quad + c_2 \exp(c_1(p - \frac{2}{L})) \cdot c \cdot 9c^2 t / L^2 \\ & \leq \frac{c^2(2c_3 + 2 + 9ct)}{L^2} c_2 \exp(c_1(p - \frac{2}{L})). \end{aligned}$$

3. For $|\mathcal{E}_3^L - \mathcal{E}_3^{L/2}|$, we first simplify

$$\mathcal{E}_3^{L/2} = \frac{2}{L} \sum_{s'=0}^s C_{t,s',0}^{L/2}(p) \Gamma_{s,s',0,0}^{L/2}(p, p),$$

and

$$\mathcal{E}_3^L = \frac{1}{L} \sum_{s'=0}^s C_{t,s',0}^L \Gamma_{s,s',0,0}^L(p, p) + C_{t,s',0}^L(p - \frac{1}{L}) \Gamma_{s,s',0,0}^L(p, p - \frac{1}{L}).$$

Again, we introduce an intermediate term

$$\tilde{\mathcal{E}}_3^L = \frac{2}{L} \sum_{s'=0}^s C_{t,s',0}^L \Gamma_{s,s',0,0}^L(p, p).$$

Then we can bound

$$\begin{aligned} & |\mathcal{E}_3^L - \mathcal{E}_3^{L/2}| \\ & \leq |\mathcal{E}_3^L - \tilde{\mathcal{E}}_3^L| + |\tilde{\mathcal{E}}_3^L - \mathcal{E}_3^{L/2}| \\ & \leq \frac{t}{L} (c_4 c_2 \exp(c_1(p - \frac{2}{L}))) / L \cdot c + c_2 \exp(c_1(p - \frac{2}{L})) \cdot c / L \\ & \quad + \frac{2t}{L} (c_5 c_2 \exp(c_1(p - \frac{1}{L}))) / L c + c_2 \exp(c_1(p - \frac{1}{L})) \cdot c / L \\ & \leq \frac{tc(c_4 + 1 + 2c_5 + 2)}{L^2} c_2 \exp(c_1(p - \frac{1}{L})). \end{aligned}$$

4. For $|\mathcal{E}_4^L - \mathcal{E}_4^{L/2}|$, we use

$$|\Gamma_{s,s',0,b}^L(p, q) - \Gamma_{s,s',0,b}^L(p - \frac{2}{L}, q) - (\Gamma_{s,s',0,b}^{L/2}(p, q) - \Gamma_{s,s',0,b}^{L/2}(p - \frac{2}{L}, q))| \leq 9c^2 t / L^2,$$

which is used in $|\mathcal{E}_2^L - \mathcal{E}_2^{L/2}|$. Finally,

$$\begin{aligned} & |\mathcal{E}_4^L - \mathcal{E}_4^{L/2}| / t \\ & \leq c_4 c_2 \exp(c_1(p - \frac{2}{L})) / L \cdot 2c / L \\ & \quad + c_2 \exp(c_1(p - \frac{2}{L})) \cdot 9c^2 t / L^2 \\ & \leq \frac{c(2c_4 + 9ct)}{L^2} c_2 \exp(c_1(p - \frac{2}{L})). \end{aligned}$$

In sum,

$$\begin{aligned} & |C_{t,s,0}^L(p + \frac{1}{L}) - C_{t,s,0}^L(p - \frac{1}{L}) - C_{t,s,0}^{L/2}(p + \frac{1}{L}) + C_{t,s,0}^{L/2}(p - \frac{1}{L})| \\ & \leq \frac{c^3 t(4ct + 2c_4 + 14) + c^2(2 + 15ct) + tc(3c_4 + 5) + c(2c_4 + 9ct)}{L^2} \cdot c_2 \exp(c_1(p - \frac{1}{L})) \\ & = \frac{c^3 t(4ct + 2c_4 + 29) + tc(3c_4 + 14) + c(2c_4 + 2c)}{L^2} \cdot c_2 \exp(c_1(p - \frac{1}{L})). \end{aligned}$$

Therefore, since $c_1 = c^3 t(4ct + 2c_4 + 29) + tc(3c_4 + 14) + c(2c_4 + 2c)$,

$$\begin{aligned} & |C_{t,s,0}^L(p + \frac{1}{L}) - C_{t,s,0}^{L/2}(p + \frac{1}{L})| \\ & \leq |C_{t,s,0}^L(p - \frac{1}{L}) - C_{t,s,0}^{L/2}(p - \frac{1}{L})| + c_1 / L^2 \cdot c_2 \exp(c_1(p - \frac{1}{L})) \\ & \leq (1 + c_1 / L) c_2 \exp(c_1(p - \frac{1}{L})) / L \\ & \leq c_2 \exp(c_1 p) / L. \end{aligned}$$

E Classification of Depthwise Parametrizations in Linear Case

We discuss the classification results on the linear residual networks with SGD training and give rigorous proofs for the claims in this simplified setting. Recall the linear residual networks:

$$\forall l \in [L], x^l = x^{l-1} + aL^{-\alpha}h^l,$$

where $h^l = W^l x^l$, and the effective learning rate of W^l is $\eta n^{-1} L^{-\gamma}$. Without loss of generality, we assume $\eta = a = 1$.

E.1 Initialization

At initialization, we have

$$\|x_0^l\rangle = \|x_0^{l-1}\rangle + L^{-\alpha}\|h_0^l\rangle,$$

where

$$\|h_0^l\rangle = \|W_0^l x_0^{l-1}\rangle = \|W_0^l x_0^{l-1}\widehat{\rangle}.$$

Since $\|x_0^{l-1}\rangle$ is independent from $\|W_0^l x_0^{l-1}\widehat{\rangle}$, we have

$$\langle x_0^l \| x_0^l \rangle = \langle x_0^{l-1} \| x_0^{l-1} \rangle + L^{-2\alpha} \langle h_0^l \| h_0^l \rangle = \langle x_0^{l-1} \| x_0^{l-1} \rangle + L^{-2\alpha} \langle x_0^{l-1} \| x_0^{l-1} \rangle = (1 + L^{-2\alpha}) \langle x_0^{l-1} \| x_0^{l-1} \rangle.$$

Using this recursion, we can write

$$\langle x_0^l \| x_0^l \rangle = (1 + L^{-2\alpha})^l \langle x_0^0 \| x_0^0 \rangle.$$

Therefore, $\langle x_0^L \| x_0^L \rangle = \Theta(1)$ iff $\alpha \geq 1/2$, otherwise $(1 + L^{-2\alpha})^L \approx e^{L^{-2\alpha+1}}$ explodes with large L .

A similar argument stands for h_0^l and f_0 . Therefore, we have proved Claim J.1.

Similarly, we can get the stability of the first backward pass, i.e., $\widetilde{\delta}x_0^l = \Theta(1)$ for $\alpha \geq 1/2$. Given $\alpha \geq 1/2$, we can also settle the size of $\widetilde{\delta}h_0$ that

$$\widetilde{\delta}h_0^l = \Theta(L^{-\alpha}),$$

which implies

$$\Delta W_1^l = L^{-\gamma+\alpha} \cdot \widetilde{\delta}h_0^l \otimes x_0^{l-1}.$$

E.2 After the first step of gradient update

Now we look at the second forward pass, and assume the input is the same, i.e., $\|x_1^0\rangle = \|x_0^0\rangle$, we have

$$\|x_1^l\rangle = \|x_1^{l-1}\rangle + L^{-\alpha}(\|W_0^l x_1^{l-1}\widehat{\rangle} + \|W_0^l x_1^{l-1}\rangle + \|\Delta W_1^l \|x_1^{l-1}\rangle)$$

where $\|\Delta W_1^l\| = -L^{-\gamma} \|\widetilde{\delta}h_0^l\rangle \langle x_0^{l-1}\| = -L^{-\gamma} \|\widetilde{\delta}x_0^l\rangle \langle x_0^{l-1}\|$, and $\|\widetilde{\delta}h_0^l\rangle \stackrel{\text{def}}{=} L^\alpha \|\widetilde{\delta}h_0^l\rangle$ is the normalized version of $\|\widetilde{\delta}h_0^l\rangle$, which happens to equal to $\|\widetilde{\delta}x_0^l\rangle$. By the definition of $\|W_0^l x_1^{l-1}\widehat{\rangle}$ and $\|W_0^l x_1^{l-1}\rangle$, we get a similar formula to the Depth- μ P case:

$$\|x_1^l\rangle = \|x_1^{l-1}\rangle + L^{-\alpha} \|W_0^l x_1^{l-1}\widehat{\rangle} + L^{-\alpha} \|\widetilde{\delta}x_0^l\rangle \left(\frac{\partial \|x_1^{l-1}\rangle}{\partial \|W_0^{l\top} \widetilde{\delta}x_0^l\rangle} - L^{-\gamma} \langle x_0^{l-1} \| x_1^{l-1} \rangle \right).$$

Now we write $b^l = L^\gamma \frac{\partial \|x_1^{l-1}\rangle}{\partial \|W_0^{l\top} \widetilde{\delta}x_0^l\rangle}$ and $c^l = -\langle x_0^{l-1} \| x_1^{l-1} \rangle$, then

$$\|x_1^l\rangle = \|x_1^{l-1}\rangle + L^{-\alpha} \|W_0^l x_1^{l-1}\widehat{\rangle} + L^{-\alpha-\gamma} (b^l + c^l) \|\widetilde{\delta}x_0^l\rangle.$$

By expanding $\|\widetilde{\delta}x_0^{l-1}\rangle = \|\widetilde{\delta}x_0^l\rangle + L^{-\alpha} \|W_0^{l\top} \widetilde{\delta}x_0^l\rangle = \|\widetilde{\delta}x_0^l\rangle + \sum_{m=l}^L L^{-\alpha} \|W_0^{m\top} \widetilde{\delta}x_0^m\rangle$, we have

$$\begin{aligned} \|x_1^l\rangle &= \|x_1^{l-1}\rangle + L^{-\alpha} \|W_0^l x_1^{l-1}\widehat{\rangle} + L^{-\alpha-\gamma} (b^l + c^l) \left(\|\widetilde{\delta}x_0^l\rangle + \sum_{m=l+1}^L L^{-\alpha} \|W_0^{m\top} \widetilde{\delta}x_0^m\rangle \right) \\ &= \|x_1^0\rangle + \sum_{m=1}^l L^{-\alpha} \|W_0^m x_1^{m-1}\widehat{\rangle} + \sum_{m=1}^l L^{-\alpha-\gamma} (b^m + c^m) \|\widetilde{\delta}x_0^m\rangle \\ &\quad + \sum_{m=2}^L L^{-\alpha-\gamma} \sum_{l'=1}^{\min\{m-1, l\}} (b^{l'} + c^{l'}) L^{-\alpha} \|W_0^{m\top} \widetilde{\delta}x_0^{m\prime}\rangle. \end{aligned} \tag{4}$$

Note the four terms in eq. (4) are independent of each other.

Now it is easy to compute c^l because only the first two terms in eq. (4) have correlation with x_0^l :

$$c^l = c^{l-1}(1 + L^{-2\alpha}) = \Theta(1)$$

with $\alpha \geq 1/2$. For b^l , we have the following recursive formula:

$$b^{l+1} = L^{-2\alpha} \sum_{m=1}^l (b^m + c^m) = \Theta(l \cdot L^{-2\alpha}).$$

Stable during training and nontrivial. Finally, we can reason about the \mathring{f}_1 (note $\mathring{f}_0 = 0$, so $\Delta \mathring{f}_1 = \mathring{f}_1$), which indicates whether the parametrization is stable during the first step¹⁰, and whether the parametrization is nontrivial for the first step:

$$\mathring{f}_1 = \langle nV \mathbb{[}x_1^L \mathbb{]} \rangle = \sum_{m=1}^L L^{-\alpha-\gamma} (b^m + c^m) \chi_0 = \Theta(L^{1-\alpha-\gamma}).$$

Therefore, we have proved Claim J.2 that the parametrization is stable during training iff $\alpha + \gamma \geq 1$, and is nontrivial iff $\alpha + \gamma \leq 1$.

Faithfulness. Although there is no activation in the linear case, we still prove Claim J.3 to enlighten the proof of the general case.

At the initialization, h_0^l and x_0^{l-1} have the same size, therefore, faithfulness is equivalent to stability, which means it happens iff $\alpha \geq 1/2$.

During training, we can expand $\mathbb{[}h_1^l \mathbb{]}$ in a similar way to eq. (4) as

$$\mathbb{[}h_1^l \mathbb{]} = \mathbb{[}W_0^l x_1^{l-1} \mathbb{]} + L^{-\gamma} (b^l + c^l) \left(\mathbb{[}\tilde{\delta} x_0^L \mathbb{]} + \sum_{m=l+1}^L L^{-\alpha} \mathbb{[}W_0^{m \top} \tilde{\delta} x_0^m \mathbb{]} \right) = \Theta(1 + L^{-\gamma}).$$

Therefore, it is faithful iff $\gamma \geq 0$. It is equivalent to $\alpha \leq 1$ because we have $\alpha + \gamma = 1$.

Feature diversity exponent. To simplify the analysis, we assume that ϵL is always an integer. We first expand $x_1^{l+\epsilon L} - x_1^l$

$$\begin{aligned} \mathbb{[}x_1^{l+\epsilon L} \mathbb{]} - \mathbb{[}x_1^l \mathbb{]} &= \sum_{m=l+1}^{l+\epsilon L} L^{-\alpha} \mathbb{[}W_0^m x_1^{m-1} \mathbb{]} + \sum_{m=l+1}^{l+\epsilon L} L^{-\alpha-\gamma} (b^m + c^m) \mathbb{[}\tilde{\delta} x_0^L \mathbb{]} \\ &\quad + \sum_{m=2}^L L^{-\alpha-\gamma} \sum_{l'=\min\{m-1, l\}+1}^{\min\{m-1, l+\epsilon L\}} (b^{l'} + c^{l'}) L^{-\alpha} \mathbb{[}W_0^{m \top} \tilde{\delta} x_0^{m'} \mathbb{]}. \end{aligned}$$

With $\alpha + \gamma = 1$, it is clear that the first term is $\Theta(L^{-\alpha} \sqrt{\epsilon L}) = \Theta(\epsilon^{1/2} L^{-\alpha+1/2})$, the second term has size $\Theta(\epsilon)$, and the third term has size $\Theta(\sqrt{L} \cdot \epsilon L^{-\alpha}) = \Theta(\epsilon L^{-\alpha+1/2})$. Therefore, there are only two cases here: if $\alpha = 1/2$, the overall size is $\Theta(\epsilon^{1/2} + \epsilon) = \Theta(\epsilon^{1/2})$; if $\alpha > 1/2$, the first and the third term vanish as $L \rightarrow \infty$, so the overall size is $\Theta(\epsilon)$. In sum, we have proved Claims J.4 and J.5.

Layerwise linearization. Claim J.6 is trivial in this simplified setting because layerwise linearization is always true for linear nets. To enlighten the proof of the general case, we recap that $\mathbb{[}\Delta W_1^l \mathbb{]} \mathbb{[}x_1^{l-1} \mathbb{]} = L^{-\gamma} c^l \mathbb{[}\tilde{\delta} x_0^l \mathbb{]} = \Theta(L^{-\gamma})$, which is much smaller than $\mathbb{[}W_0^l x_1^{l-1} \mathbb{]} = \Theta(1)$ when $\gamma > 0$. If there were an activation function, the linearization would bring an error of $o(L^{-\gamma})$ in h_1^l , which means an error of $o(L^{-\gamma-\alpha}) = o(L^{-1})$ to x_1^l .

¹⁰We need Δx and Δh for stability, but they are similar to $\Delta \mathring{f}_1$.

E.3 Beyond one step

The argument above is in general tracking the derivatives and covariance, in other words, Γ and C in the Depth- μ P case.

Now we generalize Lemma C.3, and obtain the following recursion for Γ and C

$$\begin{aligned}\Gamma_{t,r,0,b}\left(\frac{l}{L}, q\right) &= \Gamma_{t,r,0,b}\left(\frac{l-1}{L}, q\right) + L^{1/2-\alpha}\mathbb{I}_{[(t=r)\wedge(b=0)\wedge(l=\lceil Lq\rceil)]} \\ &\quad + L^{-\alpha-\gamma}\sum_{s=0}^{t-1}\Gamma_{s,r,1,b}\left(\frac{l}{L}, q\right)\left(L^{\gamma-1/2}\Gamma_{t,s,0,1}\left(\frac{l-1}{L}, \frac{l}{L}\right) - C_{t,s,0}\left(\frac{l}{L}\right)\right).\end{aligned}$$

$$\begin{aligned}\Gamma_{t,r,1,b}\left(\frac{l-1}{L}, q\right) &= \Gamma_{t,r,1,b}\left(\frac{l}{L}, q\right) + L^{1/2-\alpha}\mathbb{I}_{[(t=r)\wedge(b=1)\wedge(l=\lceil Lq\rceil)]} \\ &\quad + L^{-\alpha-\gamma}\sum_{s=0}^{t-1}\Gamma_{s,r,0,b}\left(\frac{l-1}{L}, q\right)\left(L^{\gamma-1/2}\Gamma_{t,s,1,0}\left(\frac{l}{L}, \frac{l}{L}\right) - C_{t,s,1}\left(\frac{l}{L}\right)\right).\end{aligned}$$

$$C_{t,s,a}(p) = \sum_{t'=-1}^t \sum_{s'=-1}^s \sum_{b\in\{0,1\}} \int_0^1 \Gamma_{t,t',a,b}(l/L, q) C_{t',s',b}(q) \Gamma_{s,s',a,b}(l/L, q) dq,$$

where $l = \lceil Lp \rceil - 1$ if $a = 0$, and $l = \lceil Lp \rceil$ if $a = 1$.

Then all the claims can be reasoned by tracking the order of Γ and C .

Distinguish parametrizations with $\alpha + \gamma = 1$ and $\alpha \leq 1$. The parametrizations with $\alpha + \gamma = 1$ and $\alpha \leq 1$ are all nontrivial, stable, and faithful. However, there is a large gap between $\alpha = 1/2$ (Depth- μ P) and $\alpha > 1/2$ in terms of the difficulty of tracking Γ and C . For $\alpha > 1/2$, we can see that $C_{t,s,a} = \Theta(1)$, $\Gamma_{t,-1,a,b} = \Theta(1)$ and $\Gamma_{t,s,a,b} = o(1)$ for $s \geq 0$. In this case, we can simplify the recursion by ignoring $\Gamma_{t,s,a,b}$ with $s \geq 0$:

$$\begin{aligned}\Gamma_{t,-1,0,b}\left(\frac{l}{L}\right) &\approx \Gamma_{t,-1,0,b}\left(\frac{l-1}{L}\right) - \frac{1}{L}\sum_{s=0}^{t-1}\Gamma_{s,-1,1,b}\left(\frac{l}{L}\right)C_{t,s,0}\left(\frac{l}{L}\right).\end{aligned}$$

$$\begin{aligned}\Gamma_{t,-1,1,b}\left(\frac{l-1}{L}\right) &\approx \Gamma_{t,-1,1,b}\left(\frac{l}{L}\right) - \frac{1}{L}\sum_{s=0}^{t-1}\Gamma_{s,-1,0,b}\left(\frac{l-1}{L}\right)C_{t,s,1}\left(\frac{l}{L}\right).\end{aligned}$$

$$C_{t,s,a}(p) \approx \sum_{b\in\{0,1\}} \Gamma_{t,-1,a,b}(l/L)\Gamma_{s,-1,a,b}(l/L),$$

where $l = \lceil Lp \rceil - 1$ if $a = 0$, and $l = \lceil Lp \rceil$ if $a = 1$. Note $\Gamma_{t,-1,a,b}(p, q)$ is simplified to a function that only depends on p because $\Gamma_{t,-1,a,b}(p, q)$ is constant when fixing p .

This simplification means the randomness in any W_0^l does not have an effect on the dynamics in the infinite depth limit — the complicated functional integrals for $\alpha = 1/2$ in Proposition C.4 are simplified to be ODEs when $\alpha > 1/2$. This ODE dynamic also directly implies that the feature diversity exponent is 0 for $\alpha > 1/2$.

F What Causes Hyperparameter Transfer?

In a popular misconception, hyperparameter transfer is implied by the existence of a limit. For example, the fact that μ P transfers hyperparameters, in this misconception, is because of the existence of the feature learning limit (aka the μ limit), the limit of μ P as width goes to infinity. However, this is not the case. Indeed, there are a plethora of infinite-width limits, such as the NTK limit, but there can only be one way how the optimal hyperparameters scale, so existence cannot imply transfer. In a stronger version of this misconception, transfer is implied by the existence of a “feature learning” limit. But again, this is False, because there are infinite number of feature learning limits (where the μ limit is the unique maximal one).

Instead, what is true is that the *optimal* limit implies the transfer of *optimal* hyperparameters. For example, in the width limit case, μP is the unique parametrization that yields a maximal feature learning limit. Compared to all other limits, this is obviously the optimal one. Hence μP can transfer hyperparameters across width.

So far, there is no *a priori* definition for the “optimality” of a limit: One can only tell by *classifying* all possible limits; it turns out only a small number of different behavior can occur in the limit, and thus one can manually inspect for which limit is the optimal one.

Similarly, in this work, to *derive* a depthwise scaling that allows transfer, we need to *classify* all possible infinite depth limits — and Depth- μP will turn out to be optimal in a sense that we define later in the paper.¹¹ More interestingly than the width case, here we have multiple modes of feature learning when taking the depth limit and it is important to discern which mode of feature learning is optimal. Thus, again, it is *insufficient* to derive any one limit, even with feature learning, and be able to infer it yields HP transfer.

In appendix N, we provide experiments with $1/L$ block scaling $(\alpha, \gamma) = (1, 0)$, aka ODE scaling, which provably induces feature learning in the infinite-depth limit, but is sub-optimal. Our results show a significant shift in the optimal learning rate with this parametrization.

G Notations

This section provides an introduction to the new TP notations from [24]. We only require the definition of the inner and outer products in this paper.

Averaging over n When $x \in \mathbb{R}^n$, we always use greek subscript $\alpha, \beta, \dots \in [n]$ to index its entries. Then $\langle x_\alpha \rangle_\alpha$ denotes its average entry. This notation will only be used to average over n -dimensions, but not over constant dimensions.

G.1 The Tensor Program Ansatz: Representing Vectors via Random Variables

From the Tensor Programs framework [25], we know that as width becomes large, the entries of the (pre-)activation vectors and their gradients will become roughly iid, both at initialization and training. Hence any such vector’s behavior can be tracked via a random variable that reflects the distribution of its entries. While we call this the “Tensor Program Ansatz”, it is a completely rigorous calculus.

G.1.1 Ket Notation

Concretely, if $x \in \mathbb{R}^n$ is one such vector, then we write $\llbracket x \rrbracket \in \mathbb{R}$ (called a *ket*) for such a random variable, such that x ’s entries look like iid samples from $\llbracket x \rrbracket$. For any two such vectors $x, y \in \mathbb{R}^n$, $(x_\alpha, y_\alpha) \in \mathbb{R}^2$ for each α will look like iid samples from the random vector $(\llbracket x \rrbracket, \llbracket y \rrbracket)$, such that, for example, $\lim_{n \rightarrow \infty} \frac{x^\top y}{n} = \mathbb{E} \llbracket x \rrbracket \cdot \llbracket y \rrbracket$, which we write succinctly as just $\langle x \llbracket y \rrbracket \rangle$. Here $\langle x \llbracket \cdot \rrbracket \rangle$ is called a *bra*, interpreted as a sort of “transpose” to $\llbracket x \rrbracket$. In our convention, $\llbracket x \rrbracket$ is always a random variable independent of n and x always has $\Theta(1)$ typical entry size.¹²

This notation can be generalized to the case where $\mathbf{x} \in \mathbb{R}^{n \times k}$, $\mathbf{y} \in \mathbb{R}^{n \times j}$. In this case, we can think of $\langle \mathbf{x} \llbracket \mathbf{y} \rrbracket \rangle$ as the $k \times j$ matrix given by $(\langle x_a \llbracket y_b \rrbracket \rangle)_{\substack{1 \leq a \leq k \\ 1 \leq b \leq j}}$.

Because we will often need to multiply a ket with a diagonal matrix, we introduce a shorthand:

$$\llbracket \mathbf{x} \rrbracket_{\boldsymbol{\chi}} = \llbracket \mathbf{x} \rrbracket \text{Diag}(\boldsymbol{\chi}), \quad (5)$$

if \mathbf{x} is $n \times k$ and $\boldsymbol{\chi}$ is a k -dimensional vector.

¹¹There are important nuances here that will be spelled out in an upcoming paper. For example, if the space of hyperparameters is not chosen correctly, then it could appear that no limit is *optimal* in any manner. For example, if one in (widthwise) SP, one only thinks about the 1D space of the global learning rate, then all infinite-width limits are defective — and indeed there is no hyperparameter transfer where the bigger always does better.

¹²i.e., $\|x\|^2/n = \Theta(1)$ as $n \rightarrow \infty$

G.1.2 Outer Product

Likewise, if both x and y have shape $n \times k$, the expression

$$\llbracket x \rrbracket \langle y \rrbracket \text{ represents the limit of } xy^\top \in \mathbb{R}^{n \times n}.$$

More formally, $\llbracket x \rrbracket \langle y \rrbracket$ is defined as an operator that takes a ket $\llbracket z \rrbracket \in \mathbb{R}^j$ and return the ket

$$(\llbracket x \rrbracket \langle y \rrbracket) \llbracket z \rrbracket = \llbracket x \rrbracket (\langle y \rrbracket \llbracket z \rrbracket) \in \mathbb{R}^j$$

i.e., it returns the random vector $\llbracket x \rrbracket \in \mathbb{R}^k$ multiplied by the deterministic matrix $\langle y \rrbracket \llbracket z \rrbracket \in \mathbb{R}^{k \times j}$ on the right. This corresponds to the limit of $xy^\top z/n$. Likewise, $\llbracket x \rrbracket \langle y \rrbracket$ acts on a bra $\langle w \rrbracket \in \mathbb{R}^j$ by

$$\langle w \rrbracket (\llbracket x \rrbracket \langle y \rrbracket) = (\langle w \rrbracket \llbracket x \rrbracket) \langle y \rrbracket \in \mathbb{R}^j.$$

which corresponds to the limit of $\frac{1}{n} w^\top xy^\top$. This definition of $\llbracket x \rrbracket \langle y \rrbracket$ makes the expressions

$$\llbracket x \rrbracket \langle y \rrbracket \llbracket z \rrbracket, \quad \langle w \rrbracket \llbracket x \rrbracket \langle y \rrbracket, \quad \langle w \rrbracket \llbracket x \rrbracket \langle y \rrbracket \llbracket z \rrbracket$$

unambiguous (since any way of ordering the operations give the same answer).

Remark G.1 (Potential Confusion). One should *not* interpret $\llbracket x \rrbracket \langle y \rrbracket$ as the scalar random variable $\llbracket x \rrbracket \cdot \llbracket y \rrbracket = \sum_{i=1}^k \llbracket x^i \rrbracket \llbracket y^i \rrbracket$, which would act on a ket $\llbracket z \rrbracket$ to produce $(\langle x \rrbracket \cdot \langle y \rrbracket) \llbracket z \rrbracket = \mathbb{E}(\llbracket x \rrbracket \cdot \llbracket y \rrbracket) \llbracket z \rrbracket$, which is deterministic. On the other hand, $\llbracket x \rrbracket \langle y \rrbracket \llbracket z \rrbracket$ is always a linear combination of $\llbracket x \rrbracket$, a nondeterministic random variable in general. In particular, any correlation between $\llbracket x \rrbracket$ and $\llbracket y \rrbracket$ does not directly play a role in their outer product $\llbracket x \rrbracket \langle y \rrbracket$: we always have $\llbracket x \rrbracket \langle y \rrbracket \llbracket z \rrbracket = \llbracket x \rrbracket \langle y \rrbracket^{\llbracket \square \rrbracket} \llbracket z \rrbracket^{\llbracket \square \rrbracket}$, where $(\llbracket y \rrbracket^{\llbracket \square \rrbracket}, \llbracket z \rrbracket^{\llbracket \square \rrbracket})$ is an iid copy of $(\llbracket y \rrbracket, \llbracket z \rrbracket)$ independent from $\llbracket x \rrbracket$.

Outer Product with Diagonal Inserted Finally, if $\chi \in \mathbb{R}^k$ is deterministic, then (consistent with eq. (5)) we define $\llbracket x \rrbracket_\chi \langle y \rrbracket$ as the operator that acts on kets $\llbracket z \rrbracket \in \mathbb{R}^j$ by

$$(\llbracket x \rrbracket_\chi \langle y \rrbracket) \llbracket z \rrbracket = \llbracket x \rrbracket_\chi \langle y \rrbracket \llbracket z \rrbracket = \llbracket x \rrbracket \text{Diag}(\chi) \langle y \rrbracket \llbracket z \rrbracket \in \mathbb{R}^j.$$

Morally, $\llbracket x \rrbracket_\chi \langle y \rrbracket$ is just a shorter way of writing $\llbracket x \rrbracket \text{Diag}(\chi) \langle y \rrbracket$ and represents the limit of $x \text{Diag}(\chi) y^\top$. In particular, $\llbracket x \rrbracket_1 \langle y \rrbracket = \llbracket x \rrbracket \langle y \rrbracket$.

G.1.3 Nonlinear Outer Product

If $xy^\top \in \mathbb{R}^{n \times n}$ is the (linear) outer product of two vectors $x \in \mathbb{R}^n$ and $y \in \mathbb{R}^n$, then $\phi(xy^\top)$, the entrywise application of nonlinear $\phi: \mathbb{R} \rightarrow \mathbb{R}$ to xy^\top , is a kind of *nonlinear outer product*. Passing to the ket notation, in general we define $\phi(\llbracket x \rrbracket \langle y \rrbracket)$ as the operator that acts on kets as

$$\phi(\llbracket x \rrbracket \langle y \rrbracket) \llbracket z \rrbracket \stackrel{\text{def}}{=} \mathbb{E}_{\llbracket \square \rrbracket} \phi \left(\sum_{i=1}^k \llbracket x^i \rrbracket \llbracket y^i \rrbracket^{\llbracket \square \rrbracket} \right) \llbracket z \rrbracket^{\llbracket \square \rrbracket}$$

where $(\llbracket y^1 \rrbracket^{\llbracket \square \rrbracket}, \dots, \llbracket y^k \rrbracket^{\llbracket \square \rrbracket}, \llbracket z \rrbracket^{\llbracket \square \rrbracket})$ is an iid copy of $(\llbracket y^1 \rrbracket, \dots, \llbracket y^k \rrbracket, \llbracket z \rrbracket)$ independent from $\llbracket x \rrbracket$ and the expectation is taken only over the former. This is just like, in the finite n case,

$$\phi(xy^\top) z/n = \phi \left(\sum_{i=1}^k x^i y^{i\top} \right) z/n.$$

Moreover, if $\llbracket w \rrbracket \in \mathbb{R}^j$, $\llbracket z \rrbracket \in \mathbb{R}^k$, then

$$\begin{aligned} \langle w \rrbracket \phi(\llbracket x \rrbracket \langle y \rrbracket) \llbracket z \rrbracket &= \langle w \rrbracket \phi(\llbracket x \rrbracket \langle y \rrbracket^{\llbracket \square \rrbracket}) \llbracket z \rrbracket^{\llbracket \square \rrbracket} \in \mathbb{R}^{j \times k} \\ &= \mathbb{E} \phi \left(\sum_{i=1}^k \llbracket x^i \rrbracket \llbracket y^i \rrbracket^{\llbracket \square \rrbracket} \right) (\llbracket w \rrbracket \otimes \llbracket z \rrbracket^{\llbracket \square \rrbracket}) \end{aligned}$$

where \otimes denotes outer product of vectors and expectation is taken over everything.

More generally, if $\phi: \mathbb{R}^t \rightarrow \mathbb{R}$, then $\phi(\llbracket x_1 \rrbracket \langle y_1 \rrbracket, \dots, \llbracket x_t \rrbracket \langle y_t \rrbracket)$ is an operator taking kets to kets, defined by

$$\phi(\llbracket x_1 \rrbracket \langle y_1 \rrbracket, \dots, \llbracket x_t \rrbracket \langle y_t \rrbracket) \llbracket z \rrbracket \stackrel{\text{def}}{=} \mathbb{E}_{\llbracket \square \rrbracket} \phi \left(\sum_{i=1}^k \llbracket x_1^i \rrbracket \llbracket y_1^i \rrbracket^{\llbracket \square \rrbracket}, \dots, \sum_{i=1}^k \llbracket x_t^i \rrbracket \llbracket y_t^i \rrbracket^{\llbracket \square \rrbracket} \right) \llbracket z \rrbracket^{\llbracket \square \rrbracket}$$

Remark G.2 (Potential Confusion). Note $\phi(\llbracket x \rrbracket \langle y \rrbracket)$ is not the image of the operator $\llbracket x \rrbracket \langle y \rrbracket$ under ϕ in the continuous function calculus of operators, but rather a “coordinatewise application” of ϕ . For example, if $\phi(t) = t^2$, then $\phi(\llbracket x \rrbracket \langle y \rrbracket)$ is *not* $\llbracket x \rrbracket \langle y \rrbracket \llbracket x \rrbracket \langle y \rrbracket$, the latter being what typically “squaring an operator” means, but rather $\llbracket x \rrbracket^2 \langle y \rrbracket^2 = \llbracket x \odot x \rrbracket \langle y \odot y \rrbracket$.

G.1.4 Comparison with Previous Z^\bullet Notation

For readers familiar with the *Tensor Programs* papers, this new “bra-ket” notation (aka Dirac notation) relates to the old Z^\bullet notation by

$$\llbracket x \rrbracket = Z^x, \quad \langle x \rrbracket y \rangle = \mathbb{E} Z^x Z^y.$$

The new notation’s succinctness of expectation inner product should already be apparent. Furthermore, the old notation is not very compatible with multi-vectors whereas $\llbracket x \rrbracket$ makes it clear that \rangle represents the constant dimension side. Consequently, (nonlinear) outer product is awkward to express in it, especially when its contraction with random variables requires an explicit expectation symbol \mathbb{E} .

H Infinite-Width Limit with the Bra-ket notation

As before, when the width n of the program goes to infinity, one can infer how the program behaves via a calculus of random variables. We define them below via the new ket notation instead of the earlier Z notation.

Ket Construction. We recursively define the random variable $\llbracket x \rrbracket$ (called a *ket*) for each vector x and deterministic number $\hat{\theta}$ for each scalar θ in the program. For a vector Wx in the program, we also define random variables $\llbracket Wx \hat{\cdot} \rrbracket$ and $\llbracket Wx \dot{\cdot} \rrbracket$ (called *hat-ket* and *dot-ket* respectively) such that $\llbracket Wx \dot{\cdot} \rrbracket = \llbracket Wx \hat{\cdot} \rrbracket + \llbracket Wx \dot{\cdot} \rrbracket$. These are the same as \hat{Z} and \dot{Z} in the old TP notation [25] and they satisfy

Hat All hat-kets are jointly Gaussian with zero-mean and covariance¹³

$$\text{Cov}(\llbracket Wx \hat{\cdot} \rrbracket, \llbracket Uy \hat{\cdot} \rrbracket) = \mathbb{I}(W = U) \langle x \rrbracket y \rangle \quad (6)$$

Dot Every dot-ket is a linear combination of previous kets, expressed by the following equation

$$\llbracket Wx \dot{\cdot} \rrbracket \stackrel{\text{def}}{=} \sum_{y \in \mathbf{x}} \llbracket y \rrbracket \mathbb{E} \frac{\partial \llbracket x \rrbracket}{\partial \llbracket W^\top y \hat{\cdot} \rrbracket} \quad (7)$$

eq. (7) is the same equation as in [25, Zdot] but formulated much more succinctly in the bra-ket notation:

$$[25, Zdot], \quad \dot{Z}^{Wx} = \sum_{y \in \mathbf{x}} Z^y \mathbb{E} \frac{\partial Z^x}{\partial \hat{Z}^{W^\top y}}.$$

There is an alternative notion for $\llbracket Wx \dot{\cdot} \rrbracket$ in Yang and Littwin [24] that write

$$\llbracket Wx \dot{\cdot} \rrbracket = \llbracket x \rrbracket \langle W^\top x \rrbracket.$$

This is more convenient to write as we introduce the operator view.

We can see the ket $\llbracket Wx \dot{\cdot} \rrbracket$ as the result of the action of an operator on the ket $\llbracket x \rrbracket$.

¹³In eq. (6), $\mathbb{I}(W = U)$ is the deterministic number that is 1 iff W and U are the same matrix (as symbols in the program) and 0 otherwise. This should *not* be interpreted as a random variable that is 1 precisely when W and U take the same values.

Definition H.1. Let W be an initial matrix in a Tensor Program. We define $\llbracket W \rrbracket, \widehat{\llbracket W \rrbracket}, \dot{\llbracket W \rrbracket}$ to be the linear operators on kets¹⁴ that act by

$$\begin{aligned}\widehat{\llbracket W \rrbracket}x &\stackrel{\text{def}}{=} \llbracket Wx \rrbracket \\ \dot{\llbracket W \rrbracket}x &\stackrel{\text{def}}{=} \llbracket Wx \rrbracket \\ \llbracket W \rrbracket x &\stackrel{\text{def}}{=} \widehat{\llbracket W \rrbracket}x + \dot{\llbracket W \rrbracket}x.\end{aligned}$$

Any linear operator that is equal to $\llbracket W \rrbracket$ for some initial matrix W is called an *initial operator*.

We also define the adjoint relations between the operators:

$$\begin{aligned}\widehat{\llbracket W \rrbracket}^\dagger &= \dot{\llbracket W^\top \rrbracket}, \\ \dot{\llbracket W \rrbracket}^\dagger &= \widehat{\llbracket W^\top \rrbracket}, \\ \llbracket W \rrbracket^\dagger &= \llbracket W^\top \rrbracket.\end{aligned}$$

Parameter Update In the SGD case, the parameter update of W^l is simple. With the operator notation and outer product notation, we can write

$$\llbracket W_{t+1}^l \rrbracket = \llbracket W_t^l \rrbracket - \eta \llbracket \widetilde{\delta h}_t^l \rrbracket_{\chi_t} \langle x_t^{l-1} \rrbracket.$$

In this work, Δ denotes change for one step, i.e.,

$$\llbracket \Delta W_{t+1}^l \rrbracket = -\eta \llbracket \widetilde{\delta h}_t^l \rrbracket_{\chi_t} \langle x_t^{l-1} \rrbracket;$$

$\bar{\Delta}$ denotes total change, i.e.,

$$\llbracket \bar{\Delta} W_t^l \rrbracket = -\sum_{\tau=0}^{t-1} \eta \llbracket \widetilde{\delta h}_\tau^l \rrbracket_{\chi_\tau} \langle x_\tau^{l-1} \rrbracket,$$

which we write succinctly $\llbracket \bar{\Delta} W_t^l \rrbracket = -\eta \llbracket \widetilde{\delta h}_{<t}^l \rrbracket_{\chi} \langle x_{<t}^{l-1} \rrbracket$. (Compared to Yang and Littwin [24], Δ and $\bar{\Delta}$ are changed from δ and Δ because we want to use δ for gradients instead of d , which is now used for depth differentiation).

Note in the general case,

$$\llbracket \Delta W_{t+1}^l \rrbracket = -\eta \overline{\llbracket \widetilde{\delta h}_{\leq t}^l \rrbracket_{\chi_{\leq t}} \langle x_{\leq t}^{l-1} \rrbracket}$$

where

$$\overline{\llbracket \widetilde{\delta h}_{\leq t}^l \rrbracket_{\chi_{\leq t}} \langle x_{\leq t}^{l-1} \rrbracket} \stackrel{\text{def}}{=} Q_t^l(\llbracket \widetilde{\delta h}_0^l \rrbracket_{\chi_0} \langle x_0^{l-1} \rrbracket, \dots, \llbracket \widetilde{\delta h}_t^l \rrbracket_{\chi_t} \langle x_t^{l-1} \rrbracket).$$

So

$$\llbracket \bar{\Delta} W_t^l \rrbracket = -\eta \sum_{\tau=0}^{t-1} \overline{\llbracket \widetilde{\delta h}_{\leq \tau}^l \rrbracket_{\chi_{\leq \tau}} \langle x_{\leq \tau}^{l-1} \rrbracket}. \quad (8)$$

For the rest of the paper, we write $\llbracket \bar{\Delta} W_t^l \rrbracket = -\eta \llbracket \widetilde{\delta h}_{<t}^l \rrbracket_{\chi} \langle x_{<t}^{l-1} \rrbracket$ for convenience. The generalization to eq. (8) follows Yang and Littwin [24].

I Preliminaries for the General Case

For the general case, we recall and extend the notation from the previous sections and also define new ones.

¹⁴ To be rigorous, we need to specify the ‘‘Hilbert space’’ of kets. This is somewhat pedantic and not crucial to the key points of this paper, but the Hilbert space can be constructed as follows: Let $\sigma(\pi)$ be the σ -algebra generated by the kets of the program π . Let $\Sigma(\pi) \stackrel{\text{def}}{=} \bigcup_{\pi' \supseteq \pi} \sigma(\pi')$ be the union (more precisely, the direct limit) of $\sigma(\pi')$ over all programs π' extending π . Then the Hilbert space in question is the L^2 space of random variables over the Σ of our program.

Notation Let L be the depth of the network, i.e., the number of residual blocks, and n be the width of the network, i.e. the dimension of all hidden representations x^0, \dots, x^L . Let $\xi \in \mathbb{R}^{d_{\text{in}}}$ be the input of the network, $U \in \mathbb{R}^{n \times d_{\text{in}}}$ be the input layer, and $V \in \mathbb{R}^{n \times e}$ be the output layer, so that $x^0 = U\xi$ and the model output w.r.t. ξ is $f(\xi) \triangleq V^\top x^L$. Let ℓ be the loss function absorbing the label, and δx^l be the gradient of x^l w.r.t. the loss. We denote variables at t -th training step by adding t as a subscript, e.g., the input at step t is ξ_t ¹⁵, the hidden representation of l -th layer at step t is x_t^l , and the model output at step t is f_t . Let T be the number of training steps.

I.1 Unified Scaling for SGD, Adam, and All Entrywise Optimizers

We extend the definition of entrywise update ([24]) for depth scaling, allowing us to study the unified depth scaling for SGD, Adam, and other optimization algorithms that perform only entrywise operations.

Definition I.1. A gradient-based update of parameter w with both width and depth scaling is defined by a set of functions $\mathbf{Q} = \{Q_t : \mathbb{R}^{t+1} \rightarrow \mathbb{R}\}_{t \geq 0}$, and $c, d, \delta, \gamma, \eta$. The update at time t of the optimization is

$$w \leftarrow w - \eta n^{-c} L^{-\gamma} Q_t(n^d L^\delta g_0, \dots, n^d L^\delta g_t),$$

where $g_s, s = 0, \dots, t$, are the gradients of w at time s .

For SGD, $Q_t(n^d L^\delta g_0, \dots, n^d L^\delta g_t) = n^d L^\delta g_t$, and the ‘‘true’’ learning rate is $\eta n^{-c+d} L^{-\gamma+\delta}$. For Adam,

$$Q_t(n^d L^\delta g_0, \dots, n^d L^\delta g_t) = \frac{\frac{1-\beta_1}{1-\beta_1^{t+1}} \sum_{s=0}^t \beta_1^{t-s} n^d L^\delta g_s}{\sqrt{\frac{1-\beta_2}{1-\beta_2^{t+1}} \sum_{s=0}^t \beta_2^{t-s} (n^d L^\delta g_s)^2 + \epsilon}},$$

and the ‘‘true’’ learning rate is $\eta n^{-c} L^{-\gamma}$.

The purpose of multiplying the gradients $n^d L^\delta$ before Q_t is to make sure the inputs to Q_t are $\Theta(1)$ w.r.t. n and L ¹⁶; otherwise, the update might be trivial when n and L become large. For example, if gradients are $o(1)$ entrywise, then, in Adam, directly feeding gradients to Q_t will always give an output of 0 because of the constant $\epsilon > 0$.

In this paper, we will only consider d, δ such that $n^d L^\delta g$ is $\Theta(1)$.¹⁷ As a result, the output of Q_t is also $\Theta(1)$ in general. Therefore, $n^{-c} L^{-\gamma}$ decides the scale of the update and should be our focus. We call $\eta n^{-c} L^{-\gamma}$ the *effective learning rate*.

I.2 μP and Widthwise Scaling

Maximal update parametrization (μP) [21] considers the change of initialization and learning rate of each weight matrix in the network when width scales up.¹⁸ It provides a unique initialization and learning rate of each weight matrix as a function of width n that makes the update of each weight matrix maximal (up to a constant factor). The benefit of μP is not only the theoretical guarantee but also the hyperparameter stability when scaling up the width [23].

In this paper, we assume the widthwise scaling follows μP . That is, the c in the effective learning rate $\eta n^{-c} L^{-\gamma}$ and the initialization variance of each weight matrix follows Table 1.

¹⁵Here, the input is used to perform one gradient step at training step t . We will see later that our claims should in principle hold for batched versions of the training algorithm.

¹⁶It is called faithfulness in Yang and Littwin [24].

¹⁷Note $c, d, \delta, \gamma, \eta$ in Definition I.1 can be different for parameters, so it is possible to make every parameter to satisfy the condition.

¹⁸Reparametrization is also included in the original μP , but it is not necessary for the purpose of this paper.

Table 1: Widthwise scaling of μP , where c (defined in Definition I.1) describes the widthwise scaling of the effective learning rate.

	Input weights	Output weights	Hidden weights
Init. Var.	1	n^{-2}	n^{-1}
c	0	1	1

I.3 Our Setup

We consider an L -hidden-layer residual network with biasless perceptron blocks:

$$\begin{aligned} x^0 &= U\xi, \\ \forall l \in [L], \quad x^l &= L^{-\alpha} \text{MS}(\phi(h^l)) + x^{l-1}, \quad h^l = W^l x^{l-1}, \\ f &= V^\top x^L. \end{aligned}$$

where MS refers to Mean Subtraction and is given by $\text{MS}(x) = x - \langle x, \mathbf{1} \rangle / n = Gx$ with $G = I - \mathbf{1}\mathbf{1}^\top / n$, for any $x \in \mathbb{R}^n$. The initialization and learning rate of U, V follows μP . The initialization of W^l follows μP , and the learning rate of W^l is $\eta n^{-1} L^{-\gamma}$.

Mean Subtraction (MS). The use of mean subtraction is key in our analysis. First of all, it acts in a similar fashion to LayerNorm without the normalization part. We believe that mean subtraction in LayerNorm plays a more important role than the division by the norm. This is because, with the right depth scaling, it is generally the case that features will remain stable as depth grows (no exploding/vanishing), however, the correlation will induce non-zero means, hence the need for mean subtraction to avoid exploding behavior. Secondly, subtracting the mean only decreases the dimension of the layers by 1, which is insignificant in the large width limit. Lastly, one can think of a generalized parametrization in which means of the activations are computed and added back after scaling the centered block. This will result in similar scaling patterns to the ones we propose in our setup. For the sake of simplification, we do not consider this general setup in this paper.

Definition I.2. Fix a set of update functions $\mathcal{Q} = \{Q_t : \mathbb{R}^{t+1} \rightarrow \mathbb{R}\}_{t \geq 0}$. A *depthwise parametrization* of the MLP residual network above is specified by a set of numbers $\{\alpha, \gamma, \delta\}$ such that

- (a) We independently initialize each entry of W^l from $\mathcal{N}(0, n^{-1})$
- (b) The gradients of W^l are multiplied by nL^δ before being processed by Q_t : i.e., the update at time t is

$$W^l \leftarrow W^l - \eta n^{-1} L^{-\gamma} Q_t^l(nL^\delta g_0, \dots, nL^\delta g_t) \quad (9)$$

where $g_s, s = 0, \dots, t$, are the gradients of W^l at time s and Q_t is applied entrywise.

Miscellaneous notations. For a vector x , let $[x]_i$ be its i -th coordinate. For a matrix M , let $[M]_i$ be its i -th row. Let I be the identity matrix, and $\mathbf{1}$ be the full one vector. For $m \in \mathbb{N}^+$, let $[m] = \{1, \dots, m\}$. Let \otimes be the Kronecker product.

J Classification of Depthwise Parametrizations

In this section, we provide a comprehensive description of the impact of depth parametrization on stability and update size. For this purpose, we only have two scalings to keep track of: the branch multiplier and the learning rate scaling because the initialization scale is fixed by the faithfulness property (defined below). Requiring that the features don't blow up at initialization means that the branch multipliers must be at most $\Theta(1/\sqrt{L})$. Assuming the updates are faithful (i.e., input to gradient processing functions are $\Theta(1)$ entrywise), the update size can be at most $1/L$ for the hidden layers, by an (Jacobian) operator-norm argument, but potentially much less. Naively speaking, there can be a trade-off between update size and initialization: if initialization is large, then the update may need to be small so as not to blow up the other parts of the network; likewise if the initialization is small, then the update size can be larger. But one may be surprised that a careful calculation shows that there is no trade-off: we can maximize both initialization and update size at the same time.

Before delving into the details, let us first define the notions of training routine, stability, faithfulness, and non-triviality. Hereafter, all the asymptotic notations such as \mathcal{O} , Ω and o should be understood in the limit “ $n \rightarrow \infty$, then $L \rightarrow \infty$ ”. For random variables, such notations should be understood in the sense of weak convergence (convergence in distribution). When we use the notation $x = \mathcal{O}(1)$ for some vector $x = (x_1, \dots, x_n) \in \mathbb{R}^n$, it should be understood in the sense that for all $i \in [n]$, $x_i = \mathcal{O}(1)$. Lastly, we will use bold characters (e.g. \mathbf{h} instead of h) to denote ‘batched’ versions of the quantities. This is just to emphasize that the following claims should hold for batched quantities as well.

Remark: in this section, we state the results as “claims” instead of theorems. In Appendix K.4, we provide “heuristic” proofs that can be made rigorous under non-trivial technical conditions. We also showcase the correctness of the claims by proving them rigorously in our linear setting in Appendix E. We believe this additional layer of complexity is unneeded and does not serve the purpose of this paper.

Definition J.1 (Training routine). A training routine is the package of η , \mathbf{Q} , and the input batches.

Definition J.2 (Stability). We say a parametrization is

1. *stable at initialization* if

$$\mathbf{h}_0^l, \mathbf{x}_0^l = \mathcal{O}(1), \forall l \in [L], \quad \text{and} \quad \mathbf{f}_0 = \mathcal{O}(1). \quad (10)$$

2. *stable during training* if for any training routine, any time $t \geq 0, l \in [L]$, we have

$$\Delta \mathbf{h}_t^l, \Delta \mathbf{x}_t^l = \mathcal{O}(1), \forall l \in [L], \quad \text{and} \quad \Delta \mathbf{f}_t = \mathcal{O}(1),$$

where the symbol ‘ Δ ’ refers to the change after one gradient step.

We say the parametrization is *stable* if it is stable both at initialization and during training.

Definition J.3 (Faithful). We say a parametrization is *faithful at step t* if $\mathbf{h}_t^l = \Theta(1)$ for all $l \in [L]$. We say the parametrization is *faithful* if it is faithful for all t . We also say it is *faithful at initialization* (resp. faithful during training) if this is true at $t = 0$ (resp. for $t \geq 1$).

Note faithfulness here refers to “faithfulness to ϕ ”, meaning the input to ϕ is $\Theta(1)$. This is different from the definition of faithfulness in Yang and Littwin [24], where faithfulness refers to “faithfulness to Q ” meaning the input to Q is $\Theta(1)$. “faithfulness to Q ” is already assumed in this work as mentioned in Appendix I.1.

Definition J.4 (Nontriviality). We say a parametrization is *trivial* if for every training routine and any time $t \geq 1$, $\mathbf{f}_t - \mathbf{f}_0 \xrightarrow{\text{a.s.}} 0$ in the limit “ $n \rightarrow \infty$, then $L \rightarrow \infty$ ” (i.e., the function does not evolve in the infinite-width-then-depth limit). We say the parametrization is *nontrivial* otherwise.

Definition J.5 (Feature Learning). We say a parametrization induces *feature learning* in the limit “ $n \rightarrow \infty$, then $L \rightarrow \infty$ ”, if there exist a training routine, and $t \geq 1$, and any $\lambda > 0$, we have $\Delta \mathbf{h}_t^{[\lambda L]} = \Theta(1)$.

J.1 Main Claims

We are now ready to state the main results. The next claim provides a necessary and sufficient condition under which a parametrization is stable at initialization.

Claim J.1. *A parametrization is stable at initialization iff $\alpha \geq 1/2$.*

Claim J.1 is not new and similar results were reported by Hayou et al. [7]. However, Hayou et al. [7] focuses on initialization and lacks a similar stability analysis during training. In the next result, we identify two different behaviours depending on the scaling of the learning rate.

Claim J.2. *Consider a parametrization that is stable at initialization. Then the following hold (separately from each other).*

- *It is stable during training as well iff $\alpha + \gamma \geq 1$.*
- *It is nontrivial iff $\alpha + \gamma \leq 1$.*

Therefore, it is both stable and nontrivial iff $\alpha + \gamma = 1$.

From Claim J.1 and Claim J.2, having $\alpha + \gamma = 1$ and $\alpha \geq 1/2$ is a necessary and sufficient condition for a parametrization to be stable and nontrivial throughout training. In the next result, we therefore restrict our analysis to such parametrizations and study their faithfulness.

Claim J.3. *Consider a stable and nontrivial parametrization. The following hold (separately from each other).*

- *It is faithful at initialization iff $\alpha \geq 1/2$. As a result, $\alpha = 1/2$ is the minimal choice of α that guarantees faithfulness.*
- *It is faithful during training iff $\alpha \leq 1$.*

Therefore, a stable and nontrivial parametrization is faithful iff $\alpha \in [1/2, 1]$.

The first claim follows from well-known calculations of randomly initialized residual networks [7]. For the second claim, the intuition here is just that if $\alpha + \gamma = 1$ and $\alpha > 1$ then $\gamma < 0$, i.e., the update size blows up with depth. This would then cause the input to the nonlinearities to blow up with size.

One might argue that faithfulness at initialization is not important (e.g. features at initialization could converge to zero without any stability or triviality issues) and what matters is faithfulness throughout training. It turns out that faithfulness at initialization plays a crucial role in the optimal use of network capacity. To see this, we first define the notion of feature diversity exponent, which relates to the similarity in the features of adjacent layers.

Definition J.6 (Feature Diversity Exponent). We say a parametrization has feature diversity exponent $\kappa \geq 0$ if κ is the maximal value such that for all $\lambda \in [0, 1]$ and sufficiently small $\epsilon > 0$, and all time t ,

$$\frac{1}{\sqrt{n}} \left\| \mathbf{x}_t^{\lfloor (\lambda + \epsilon)L \rfloor} - \mathbf{x}_t^{\lfloor \lambda L \rfloor} \right\| = \Omega(\epsilon^{1 - \kappa}),$$

where $\Omega(1)$ should be interpreted in the limit “ $n \rightarrow \infty$, then $L \rightarrow \infty$, then $\epsilon \rightarrow 0$ ”. We say a parametrization is *redundant* if $\kappa = 0$.

In other words, the feature diversity exponent κ is a measure of how different the outputs are in layers that are close to each other. With $\kappa = 0$, the output of each layer is essentially the same as the output of the previous layer in the sense that the rate of change from one layer to the next is bounded (at least locally), and hence the network is intuitively “wasting” parameters.

Claim J.4. *Consider a stable and nontrivial parametrization that is furthermore faithful during training (but not necessarily at initialization). Then it is redundant if $\alpha \in (1/2, 1]$.*

To understand the intuition behind Claim J.4, let us see what happens when $\alpha > 1/2$. In this case, the randomness of the initialization weights will have no impact on training trajectory as depth increases. To see this, consider some layer index $\lfloor \lambda L \rfloor$. The blocks are divided by L^α which is larger than the magnitude of accumulated randomness (of order $(\lambda L)^{1/2}$). This basically destroys all the randomness from initialization and therefore the randomness in the learned features will consist only of that coming from U and V (input and output matrices). When depth goes to infinity, the contribution of the randomness in two adjacent layers becomes less important, we end up with adjacent layers becoming very similar because the gradients to these layers are highly correlated.

In contrast, we have the following result, which defines Depth- μ P.

Claim J.5 (Depth- μ P). $\alpha = \gamma = 1/2$ is the unique parametrization that is stable, nontrivial, faithful, induces feature learning, and achieves maximal feature diversity with $\kappa = 1/2$.

In terms of feature diversity, a phase transition phenomenon occurs when $\alpha = 1/2$. More precisely, for Depth- μ P, we can show that $n^{-1/2} \left\| \mathbf{x}_t^{\lfloor (\lambda + \epsilon)L \rfloor} - \mathbf{x}_t^{\lfloor \lambda L \rfloor} \right\| = \mathcal{O}(\epsilon^{1/2})$ while the same quantity is $\mathcal{O}(\epsilon)$ for all $\alpha \in (1/2, 1]$, which suggests that Depth- μ P yields *rough* path for \mathbf{x}_t . This allows the features to change significantly from one layer to the next, hence efficiently using the parameters. For readers who are familiar with rough path theory, the $1/2$ continuity exponent is a result of Brownian increments in the path.¹⁹

¹⁹The reader might ask whether we can obtain an exponent smaller than $1/2$. This is indeed possible, but it will entail using correlated weights. We leave this question for future work.

Moreover, with $\alpha = 1$, there is a phenomenon of feature collapse in the sense that the features will be contained in the σ -algebra generated by the input and output layers, but contains no randomness from the hidden layers (see Appendix K.2). Intuitively, the case of $\alpha = 1$ is analogous to width situation, where deep mean field collapses to a single neuron (all neurons become essentially the same). For depth, the features (layers) are still relatively different but the redundancy does not allow significant variety in these features.

J.2 Subtlety: Layerwise (local) linearization but not global linearization

Definition J.7. We say a parametrization induces layerwise linearization iff each layer can be linearized without changing the network output when $L \rightarrow \infty$, that is, $\forall l \in [L]$,

$$L^{-\alpha} G (\phi(W_t^l \mathbf{x}_t^{l-1}) - \phi(W_0^l \mathbf{x}_t^{l-1}) - \phi'(W_0^l \mathbf{x}_t^{l-1}) \odot ((W_t^l - W_0^l) \mathbf{x}_t^{l-1})) = o(L^{-1})$$

Claim J.6. A stable and nontrivial parametrization induces layerwise linearization iff $\alpha \in [1/2, 1)$.

However, note that this does not imply the entire network is linearized (w.r.t. all the parameters in the sense of Neural Tangent Kernel). In our setup, where the input and output layers are initialized at a constant scale (w.r.t. L), it is actually not possible to have a kernel limit. Even in our linear case in Appendix C, one can see the learned model is not linear.

If the initialization of the output layer is L times larger than our setup (assuming $L \ll n$ so the widthwise scaling still follows μP), it may induce a parametrization that can linearize the entire network. In that situation, the learning rate has to be L times smaller than Depth- μP to obtain stability during training, so the change of parameters is also L times smaller, which can lead to the linearization of the entire network. Since we focus on maximal feature learning, the rigorous argument is beyond the scope of this paper.

K Heuristics for the proofs in the general case

The notation in this section is mostly defined in appendix G. The complete notation is defined in [24].

K.1 Depth- μP

Let $\text{MS}(x) = x - \langle x, 1 \rangle / n = Gx$ where $G = I - 11^\top / n$, where $x \in \mathbb{R}^n$. Recall the definition of the network and the normalized gradients

$$\begin{aligned} x^1 &= U\xi \\ h^l &= W^l x^{l-1} \\ x^l &= x^{l-1} + \frac{1}{\sqrt{L}} G \phi(h^l) \\ f(\xi) &= V^\top x^L \\ \tilde{\delta} x^L &= nV \\ \tilde{\delta} h^l &= \phi'(h^l) \odot (G \tilde{\delta} x^l) \\ \tilde{\delta} x^{l-1} &= \tilde{\delta} x^l + \frac{1}{\sqrt{L}} W^{l\top} \tilde{\delta} h^l \end{aligned}$$

where $V = \Theta(1/n)$ coordinatewise, $\tilde{\delta} x^l = \Theta(1)$ coordinatewise and $W^l = \Theta(\frac{1}{\sqrt{n}})$ coordinate-wise.

We also abuse the notation of G and use it as an operator on kets: $G|x\rangle \stackrel{\text{def}}{=} |x\rangle - \mathbb{E}|x\rangle$.

Forward. Similar to the linear case, one can show that under technical conditions (mostly on the activation function) that the infinite-depth limit of the TP follows the dynamics

$$\begin{aligned}
d\|x_t^\lambda\rangle &= \sqrt{d\lambda}G\phi\left(\|W_0^\lambda\|x_t^\lambda\rangle + \sqrt{d\lambda}\|\widetilde{\Delta W}_t^\lambda\|x_t^\lambda\rangle\right) \\
&= \sqrt{d\lambda}G\phi\left(\|W_0^\lambda\|x_t^\lambda\rangle\right) + d\lambda G\phi'\left(\|W_0^\lambda\|x_t^\lambda\rangle\right)\|\widetilde{\Delta W}_t^\lambda\|x_t^\lambda\rangle \\
&= \sqrt{d\lambda}G\phi\left(\widehat{\|W_0^\lambda\|}x_t^\lambda\rangle + \dot{\|W_0^\lambda\|}x_t^\lambda\rangle\right) + d\lambda G\phi'\left(\|W_0^\lambda\|x_t^\lambda\rangle\right)\|\widetilde{\Delta W}_t^\lambda\|x_t^\lambda\rangle \\
&= \sqrt{d\lambda}G\phi\left(\widehat{\|W_0^\lambda\|}x_t^\lambda\rangle\right) + d\lambda G\phi'\left(\|W_0^\lambda\|x_t^\lambda\rangle\right)\left(\dot{\|W_0^\lambda\|}x_t^\lambda\rangle + \|\widetilde{\Delta W}_t^\lambda\|x_t^\lambda\rangle\right)
\end{aligned}$$

where $\lambda \in [0, 1]$ refers to the fractional layer index (λ represents layer index $\lfloor \lambda L \rfloor$ as $L \rightarrow \infty$), t refers to the training step, $\|W_0^\lambda\|$ the matrix operator (defined in Appendix H), and the tilde symbol refers to the ‘‘normalized’’ version of the object, i.e., multiply the ket with $(d\lambda)^c$ for some c such that the multiplication (normalized ket) is $\Theta(1)$ w.r.t. L , and same for the normalized operators. We also simplify $\tilde{\delta}$ to δ if it is already under wider tilde symbol. The first term represents a Gaussian noise.

In the linear case, we have

$$d\|x_t^\lambda\rangle = \sqrt{d\lambda}\left(\widehat{\|W_0^\lambda\|}x_t^\lambda\rangle\right) + d\lambda\left(\dot{\|W_0^\lambda\|}x_t^\lambda\rangle + \|\widetilde{\Delta W}_t^\lambda\|x_t^\lambda\rangle\right)$$

Note

$$\dot{\|W_0^\lambda\|}x_t^\lambda\rangle = \sqrt{d\lambda}\sum_{s=0}^{t-1}\|\widetilde{\delta h}_s^\lambda\rangle\langle\nabla_{W_0^{\lambda\top}\widetilde{\delta h}_s^\lambda}\|x_t^\lambda\rangle = \sqrt{d\lambda}\|\widetilde{\delta h}_{<t}^\lambda\rangle\langle W_0^{\lambda\top}\widetilde{\delta h}_{<t}^\lambda\rangle\|x_t^\lambda\rangle$$

Using multi-vector notation, we write

$$\begin{aligned}
\|\widetilde{\Delta W}_t^\lambda\|x_t^\lambda\rangle &= -\eta\|\widetilde{\delta h}_{<t}^\lambda\rangle_{\mathbf{x}}\langle\mathbf{x}_{<t}^\lambda\|x_t^\lambda\rangle = -\eta\sum_{s<t}\|\widetilde{\delta h}_s^\lambda\rangle_{\mathbf{x}_s}\langle x_s^\lambda\|x_t^\lambda\rangle \\
\|\Delta W_t^\lambda\|x_t^\lambda\rangle &= -\eta\|\widetilde{\delta h}_{<t}^\lambda\rangle_{\mathbf{x}}\langle\mathbf{x}_{<t}^\lambda\|x_t^\lambda\rangle = -\eta\sum_{s<t}\|\widetilde{\delta h}_s^\lambda\rangle_{\mathbf{x}_s}\langle x_s^\lambda\|x_t^\lambda\rangle
\end{aligned}$$

Backward. Similar to the forward prop, we obtain the following dynamics for the infinite-depth TP

$$\begin{aligned}
-d\|\widetilde{\delta x}_\tau^\lambda\rangle &= \sqrt{d\lambda}\|W_\tau^{\lambda\top}\|\phi'(W_\tau^\lambda x_\tau^\lambda)\odot(G\widetilde{\delta x}_\tau^\lambda) \\
&= \sqrt{d\lambda}\left(\widehat{\|W_0^{\lambda\top}\|} + \dot{\|W_0^{\lambda\top}\|} + \sqrt{d\lambda}\|\widetilde{\Delta W}_\tau^\lambda\|^\dagger\right)\left[\phi'\left(\|W_0^\lambda x_\tau^\lambda\rangle + \sqrt{d\lambda}\|\widetilde{W}_0^\lambda x_\tau^\lambda\rangle + \sqrt{d\lambda}\|\widetilde{\Delta W}_\tau^\lambda\|x_\tau^\lambda\rangle\right)\|G\widetilde{\delta x}_\tau^\lambda\rangle\right] \\
&= \sqrt{d\lambda}\widehat{\|W_0^{\lambda\top}\|}\left[\phi'(\|W_0^\lambda x_\tau^\lambda\rangle)\|G\widetilde{\delta x}_\tau^\lambda\rangle\right] + \sqrt{d\lambda}\dot{\|W_0^{\lambda\top}\|}\left[\phi'(\|W_0^\lambda x_\tau^\lambda\rangle)\|G\widetilde{\delta x}_\tau^\lambda\rangle\right] + d\lambda\|\widetilde{\Delta W}_\tau^\lambda\|^\dagger\left[\phi'(\|W_0^\lambda x_\tau^\lambda\rangle)\|G\widetilde{\delta x}_\tau^\lambda\rangle\right] \\
&\quad + d\lambda\dot{\|W_0^{\lambda\top}\|}\left[\phi''(\|W_0^\lambda x_\tau^\lambda\rangle)\left\{\|\widetilde{W}_0^\lambda x_\tau^\lambda\rangle + \|\widetilde{\Delta W}_\tau^\lambda\|x_\tau^\lambda\rangle\right\}\|G\widetilde{\delta x}_\tau^\lambda\rangle\right]
\end{aligned}$$

Here the $(d\lambda)^{3/2}$ term got dropped. The individual terms can be simplified as follows

$$\begin{aligned}
\|\widetilde{\Delta W}_\tau^\lambda\|^\dagger\left[\phi'(\|W_0^\lambda x_\tau^\lambda\rangle)\|G\widetilde{\delta x}_\tau^\lambda\rangle\right] &= -\eta\|\mathbf{x}_{<\tau}^\lambda\rangle_{\mathbf{x}}\langle\widetilde{\delta h}_{<\tau}^\lambda\|\left[\phi'(\|W_0^\lambda x_\tau^\lambda\rangle)\|G\widetilde{\delta x}_\tau^\lambda\rangle\right] \approx -\eta\|\mathbf{x}_{<\tau}^\lambda\rangle_{\mathbf{x}}\langle\widetilde{\delta h}_{<\tau}^\lambda\|\widetilde{\delta h}_\tau^\lambda\rangle \\
\dot{\|W_0^{\lambda\top}\|}\left[\phi'(\|W_0^\lambda x_\tau^\lambda\rangle)\|G\widetilde{\delta x}_\tau^\lambda\rangle\right] &= \left[\|\mathbf{x}_{<\tau}^\lambda\rangle\langle W_0^\lambda \mathbf{x}_{<\tau}^\lambda\| + \|\mathbf{x}_\tau^\lambda\rangle\langle W_0^\lambda \mathbf{x}_\tau^\lambda\| \right]\left[\phi'(\|W_0^\lambda x_\tau^\lambda\rangle)\|G\widetilde{\delta x}_\tau^\lambda\rangle\right] \\
&= \|\mathbf{x}_{<\tau}^\lambda\rangle_{\mathbf{x}}\mathbb{E}\left[\phi'(\|W_0^\lambda x_\tau^\lambda\rangle)\frac{\partial\|G\widetilde{\delta x}_\tau^\lambda\rangle}{\partial\|W_0^\lambda \mathbf{x}_{<\tau}^\lambda\rangle}\right] + \|\mathbf{x}_\tau^\lambda\rangle_{\mathbf{x}}\mathbb{E}\left[\phi''(\|W_0^\lambda x_\tau^\lambda\rangle)\|G\widetilde{\delta x}_\tau^\lambda\rangle\right] \\
&= \Theta(\sqrt{d\lambda})
\end{aligned}$$

where the other terms from the product rule drops out because

$$\frac{\partial\phi'(\|W_0^\lambda x_\tau^\lambda\rangle)}{\partial\|W_0^\lambda \mathbf{x}_{<\tau}^\lambda\rangle} = \frac{\partial\|G\widetilde{\delta x}_\tau^\lambda\rangle}{\partial\|W_0^\lambda \mathbf{x}_\tau^\lambda\rangle} = 0$$

K.2 $1/L$ branches

K.2.1 Forward:

$$\begin{aligned} d\|x_t^\lambda\rangle &= d\lambda G \mathbb{E} \left[\phi \left(\|W_0^\lambda\|x_t^\lambda\rangle + \|\bar{\Delta}W_t^\lambda\|x_t^\lambda\rangle \right) \mid \|U_0, V_0\rangle \right] \\ &= d\lambda G \mathbb{E} \left[\phi \left(\widehat{\|W_0^\lambda\|}x_t^\lambda\rangle + \|\bar{\Delta}W_t^\lambda\|x_t^\lambda\rangle \right) \mid \|U_0, V_0\rangle \right] \end{aligned}$$

where the equality follows because $\|x_t^\lambda\rangle$ is contained the σ -algebra of $\|U_0, V_0\rangle$, so $\dot{\|W_0^\lambda\|}x_t^\lambda\rangle = 0$. Since $\|\bar{\Delta}W_t^\lambda\| \in \sigma(\|U_0, V_0\rangle) \otimes \sigma(\|U_0, V_0\rangle)$, $\|\bar{\Delta}W_t^\lambda\|x_t^\lambda\rangle \in \sigma(\|U_0, V_0\rangle)$, and the expectation is really just over $\widehat{\|W_0^\lambda\|}x_t^\lambda\rangle$.

K.2.2 Backward

$$\begin{aligned} -d\|\tilde{\delta}x_\tau^\lambda\rangle &= d\lambda \mathbb{E} \left[\|W_\tau^{\lambda\top}\| \phi'(W_\tau^\lambda x_\tau^\lambda) \odot (G\tilde{\delta}x_\tau^\lambda) \mid \|U_0, V_0\rangle \right] \\ &= d\lambda \mathbb{E} \left[\|\bar{\Delta}W_\tau^{\lambda\top}\| \phi'(W_\tau^\lambda x_\tau^\lambda) \odot (G\tilde{\delta}x_\tau^\lambda) \mid \|U_0, V_0\rangle \right] \end{aligned}$$

Here the $\widehat{\|W_0^{\lambda\top}\|}$ and $\dot{\|W_0^{\lambda\top}\|}$ drop out because the former is zero-mean and independent from $\|U_0, V_0\rangle$ and the latter drops out because $\|x_t^\lambda\rangle$ is contained the σ -algebra of $\|U_0, V_0\rangle$.

K.3 $1/L^\alpha$ branches, $\alpha \in (1/2, 1]$

K.3.1 Forward

$$\begin{aligned} d\|x_t^\lambda\rangle &= (d\lambda)^\alpha G \mathbb{E} \left[\phi \left(\|W_0^\lambda\|x_t^\lambda\rangle + (d\lambda)^{1-\alpha} \|\widetilde{\bar{\Delta}W_t^\lambda}\|x_t^\lambda\rangle \right) \mid \|U_0, V_0\rangle \right] \\ &= d\lambda \mathbb{E} \left[\phi' \left(\widehat{\|W_0^\lambda\|}x_t^\lambda\rangle \right) \right] G \|\bar{\Delta}W_t^\lambda\|x_t^\lambda\rangle \end{aligned}$$

because the same reason as above.

K.3.2 Backward

$$\begin{aligned} -d\|\tilde{\delta}x_\tau^\lambda\rangle &= d\lambda \mathbb{E} \left[\|W_\tau^{\lambda\top}\| \phi'(W_\tau^\lambda x_\tau^\lambda) \odot (G\tilde{\delta}x_\tau^\lambda) \mid \|U_0, V_0\rangle \right] \\ &= d\lambda \mathbb{E} \left[\|\bar{\Delta}W_\tau^{\lambda\top}\| \phi'(W_\tau^\lambda x_\tau^\lambda) \odot (G\tilde{\delta}x_\tau^\lambda) \mid \|U_0, V_0\rangle \right] \\ &= d\lambda \mathbb{E} \left[\phi'(\widehat{\|W_0^\lambda\|}x_\tau^\lambda) \right] \|\bar{\Delta}W_\tau^{\lambda\top}\| G\tilde{\delta}x_\tau^\lambda \end{aligned}$$

Here,

$$\|\phi'(W_\tau^\lambda x_\tau^\lambda) \odot (G\tilde{\delta}x_\tau^\lambda)\rangle \equiv \mathbb{E} \left[\phi'(\widehat{\|W_0^\lambda\|}x_\tau^\lambda) \right] G\|\tilde{\delta}x_\tau^\lambda\rangle + (d\lambda)^{1-\alpha} \mathbb{E} \left[\phi''(\widehat{\|W_0^\lambda\|}x_\tau^\lambda) \right] G\|\bar{\Delta}W_\tau^{\lambda\top}\|\tilde{\delta}x_\tau^\lambda\rangle$$

K.4 Justifications of the claims

Claim J.2. Stability during training when $\alpha + \gamma \geq 1$ is straightforward (some technical conditions on the activation function are required). This is because the weight updates are of order $L^{-\alpha-\gamma}$ and feature updates involve no more than L terms of size $L^{-\alpha-\gamma}$ (plus higher order terms that do not contribute to the update in the large depth limit). When $\alpha + \gamma > 1$, the contribution of a sum of at most L terms of order $L^{-\alpha-\gamma}$ will decrease to zero, and the network output f_t will converge to f_0 in this case, yielding a trivial limit. However, when $\alpha + \gamma = 1$, the updates remain important in the infinite depth limit, yielding a non-trivial limit.

Claim J.3. Consider a stable and nontrivial parametrization (i.e. $\alpha + \gamma = 1$). Faithfulness at initialization is achieved only when $\alpha \geq 1/2$. This was proven in [7] in a more general setup.

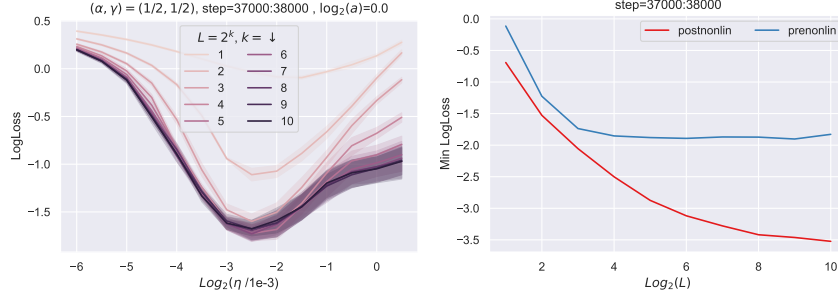


Figure 3: **Pre-Nonlin Leads to Poor Performance** Although Depth- μ P for prenonlin resnet indeed transfers hyperparameters (Left), depth gives no performance gains beyond 8 layers and the performance is dramatically worse than the post-nonlinearity resnet (Right). In right plot, the "Min LogLoss" is minimal log loss over all block multiplier and learning rate. CIFAR10+Adam. See Figure 11 for more details about the setup.

Faithfulness during training is ensured as long as $\alpha \leq 1$ because feature updates are always $\Theta(1)$ in depth. With $\alpha > 1$, $\gamma < 0$ and the weight updates explode with depth in this case, which yield exploding behaviour for h .

Claim J.4 When $\alpha \in (1/2, 1]$, we obtain smooth limiting dynamics when $L \rightarrow \infty$ as demonstrated in Appendix K.3. This limiting process is a smooth process (no Brownian jumps) that satisfies the required definition of redundancy.

Claim J.5. It remains to prove that Depth- μ P is non-redundant. This is a result of the limiting dynamics in this case (Appendix K.1). With Depth- μ P, the randomness of the initialization in the hidden layer remains present throughout training, inducing a Brownian-like term that breaks redundancy.

Claim J.6. In Depth- μ P, $W_t^l - W_0^l$ is $\Theta(1/\sqrt{L})$ which is much smaller than W_0^l . Therefore, $\phi(W_t^l x_t^{l-1}) - \phi(W_0^l x_t^{l-1}) - \phi'(W_0^l x_t^{l-1}) \odot ((W_t^l - W_0^l) x_t^{l-1}) = o(1/\sqrt{L})$, thus satisfies Definition J.7. Similar to the depth- μ P case, for $\alpha \in [1/2, 1)$, the activation in the forward pass can be linearized which indicates layerwise linearization when $\alpha + \gamma = 1$.

L Feature Diversity

In this section, we show that the choice of nonlinearity and placement of nonlinearities can affect feature diversity greatly.

L.1 Gradient Diversity

Gradient diversity is an important factor toward feature diversity. Observe that the gradient δx^l at x^l is continuous in l in the limit $L \rightarrow \infty$. In a linear model (or the pre-nonlinearity model, where nonlinearity is put before the weights), this causes $\delta h^l = L^{-\alpha} \delta x^l$ to be very similar between neighboring blocks. As a result (because the weights W^l receives an update proportional to $\delta h^l \otimes x^{l-1}$), in the next forward pass, neighboring blocks contribute very similarly to the main branch x^l . This leads to a waste of model capacity.

L.2 Pre-Nonlin Leads to Poor Performance

For example, in Figure 3, for a relu pre-nonlinearity resnet (i.e. blocks are given by $W^l \phi(x^{l-1})$ instead of $\phi(W^l x^{l-1})$), we see that although Depth- μ P indeed transfers hyperparameters (as predicted by our theory), the performance is dramatically worse than the post-nonlinearity resnet in Figure 11, and depth gives no performance gains beyond 8 layers. Specifically, it is because $\delta h^l = L^{-\alpha} \delta x^l$ like the linear case, and $\phi(x^{l-1})$ is also similar between neighboring blocks. As a result, the gradient of the weights W^l , proportional to $\delta h^l \otimes \phi(x^{l-1})$, has little diversity compared to nearby blocks.

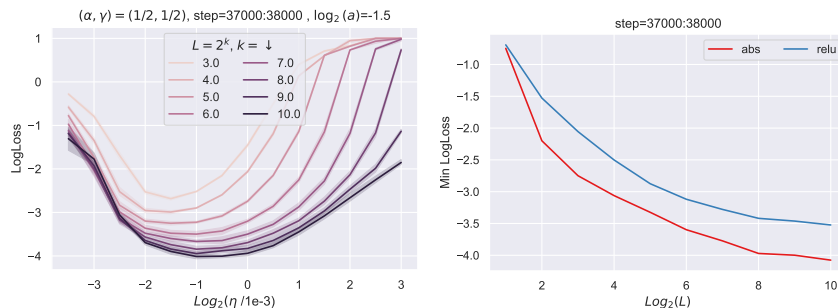


Figure 4: **Improving performance with absolute value non-linearity**, which maximizes feature diversity. (CIFAR10+Adam). See Figure 11 for more details about the setup.

L.3 Maximizing Feature Diversity with Absolute Value Nonlinearity

In a nonlinear model, we have $\delta h^l = \delta x^l \odot \phi'(h^l)$. Because h^l is almost independent from all other h^m , $m \neq l$ in the Depth- μ P limit, $\phi'(h^l)$ can serve to decorrelate the δh^l , depending on what ϕ is. For example, if ϕ is relu, then ϕ' is the step function. h^l is approximately a zero-mean Gaussian in the Depth μ P limit, so that $\phi'(h^l)$ is approximately 0 or 1 with half probability each. This decorrelates δh^l much better than the linear case. But of course, this line of reasoning naturally leads to the conclusion that $\phi' = \text{sign}$ would be the best decorrelator of δh^l and the maximizer of feature diversity (with ϕ among the class of positively 1-homogeneous functions) — then δh^l and δh^m are completely decorrelated for $l \neq m$.

Indeed, as shown in Figure 4, swapping in absolute value for ϕ dramatically improves the training performance of deep (block depth 1) resnets.

In general, in lieu of absolute value, any even nonlinearity would suffice.

L.4 Feature Diversity is in Tension with Layerwise Linearization

The reason that $\phi'(h^l)$ can decorrelate δh^l is very much related to layerwise linearization. Recall that in Depth- μ P, h^l can be decomposed to a zero-mean Gaussian part \widehat{h}^l of size $\Theta(1)$ and a correction term \dot{h}^l of size $\Theta(L^{-1/2})$ (corresponding to the decomposition $\|h^l\rangle = \|\widehat{h}^l\rangle + \|\dot{h}^l\rangle$). \widehat{h}^l is independent from \widehat{h}^m for $m \neq l$ but \dot{h}^l can be very strongly correlated to all other \dot{h}^m . Thus, $\phi'(h^l)$ can decorrelate δh^l precisely because \widehat{h}^l dominates \dot{h}^l , and this is also precisely the reason we have layerwise linearization.

In the $1/L$ scaling $(\alpha, \gamma) = (1, 0)$, \widehat{h}^l is on the same order as \dot{h}^l and layerwise linearization does not occur, but also $\phi'(h^l)$ can no longer effectively decorrelate δh^l .

Once again, we remind the reader that layerwise linearization in this case is not detrimental (in this block depth 1 case) because \widehat{h}^l in fact accumulate contributions from the learned features of all previous blocks and thus strongly depends on the learning trajectory (in contrast to the (widthwise) NTK case where \widehat{h}^l is already determined at initialization).

M Block Depth 2 and Above

Remark on notation: Here and in the next section, all big-O notation is in L only; the scaling in width is assumed to be in μ P.

In most of this work, we have considered depth-1 MLP for g^l in eq. (1), it's straightforward to derive and classify the infinite-width-then-infinite-depth limits for larger depths in each block. In particular, the following $1/\sqrt{L}$ scaling still makes sense in this more general setting with block depth k and leads to a well defined limit:

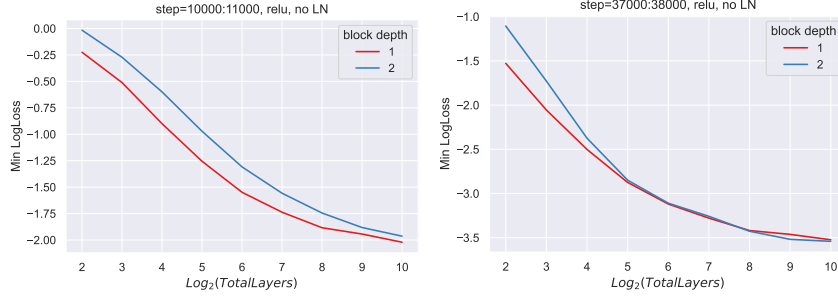


Figure 5: **Block Depth 2 < Block Depth 1, Relu.** In relu resnet with no LN, block depth 2 does worse than block depth 1 when matching total number of layers (and thus parameter count). However, training longer (38000 steps, Right) helps it catch up (compared to 11000 steps, Left). The y-axis is minimal log loss over all block multiplier and learning rate

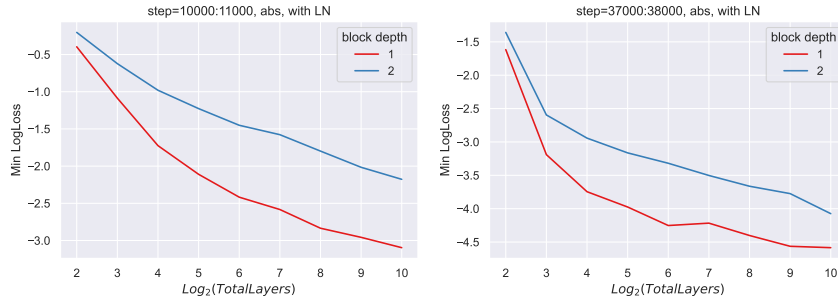


Figure 6: **Block Depth 2 < Block Depth 1, Abs.** In abs resnet with LN, block depth 2 does significantly worse than block depth 1 when matching total number of layers (and thus parameter count). Training longer (38000 steps, Right) does not close the performance gap (compared to 11000 steps, Left). The y-axis is minimal log loss over all block multiplier and learning rate

$$x^l = x^{l-1} + \frac{a}{\sqrt{L}} \cdot g^l(x^{l-1}; W^{l1}, \dots, W^{lk}), \quad \Theta(1) \text{ initialization scale, } \Theta(1/\sqrt{L}) \text{ learning rate} \tag{11}$$

This is what we call Depth- μ P in the block depth 1 case, but we shall not use this name in the general block depth case because *this parametrization is no longer optimal*.²⁰

M.1 Block Depth ≥ 2 is Defective

A very clear symptom of this is that the *performance of block-depth-2 resnets is worse than that of block-depth-1 networks*, when matching parameter count, although they can (but not always) catch up after training for a long time (figs. 5 and 6). Simultaneously, we are seeing nontrivial or even significant hyperparameter shifts as the total number of blocks increases (fig. 7).

M.2 Defect of $1/\sqrt{L}$ Scaling in Block Depth 2

The reason that the $1/\sqrt{L}$ scaling is no longer fine in the block depth ≥ 2 case is the *linearization of the multiplicative interaction* between the layers in the block. Indeed, just like the block depth 1 case, the $1/\sqrt{L}$ scaling forces the weight updates ΔW of each weight matrix to be $\Theta(\sqrt{L})$ smaller than the initialization W_0 . Thus, within the block, the training dynamics when depth L is large is in the kernel regime, where the contribution to the block output $g(x; W^\bullet)$ is only a *summation*, instead of *product*, of individual contributions from each layer’s weights updates.

²⁰What we exactly mean by *optimal* will be explained below.

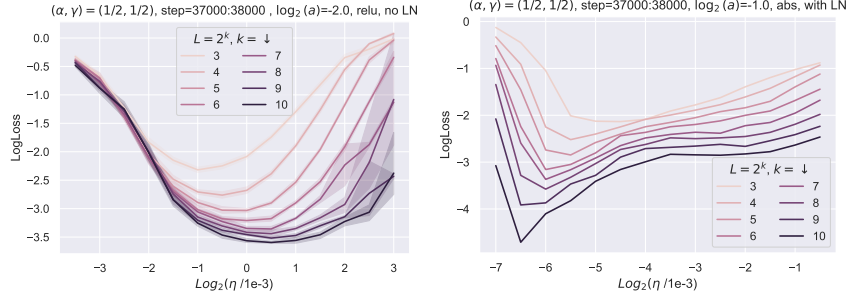


Figure 7: **Block Depth 2 Hyperparameter Shift** in relu resnet with no LN (Left) and abs resnet with LN (Right).

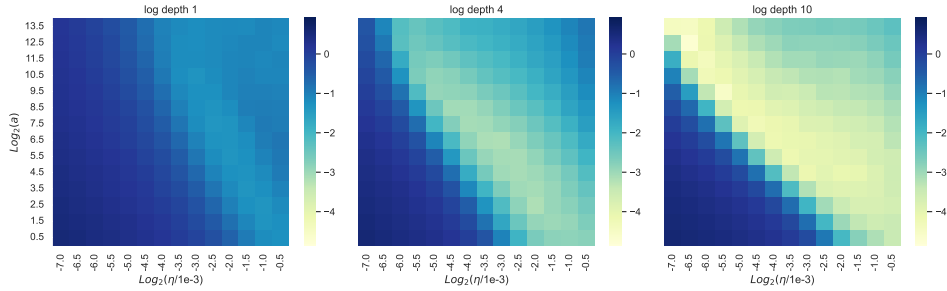


Figure 8: The "slope" of the optimal sublevel set in the (learning rate, block multiplier) space changes from -2 to -1 as depth goes from 2^1 to 2^{10} . Absolute value nonlinearity with layer normalization, block depth 2, 38000 steps of Adam on CIFAR10.

When aggregated over all L blocks, the result is that there is only multiplicative interaction across blocks but not within layers. In other words, the network output is dominated by the contributions of the form $W^{L i_L} \dots W^{1 i_1}$ where each $i_l \in [k]$ is a layer choice within block l . All other contributions (which all involve within-block interactions) are subleading.

When block depth $k = 1$ (our main subject of study in this work), *all* interactions are included but this is no longer true when $k > 1$.

In fig. 8, the heatmap of loss as a function of block multiplier and learning rate demonstrates this vividly for block depth 2.

Small depth The optimal sublevel set of (learning rate, block multiplier) has slope ≈ -2 when the number of blocks is 2^1 . In other words, around the optimum, double the learning rate while dividing the block multiplier by 4 has similar performance. This is because ΔW^{l1} and ΔW^{l2} interact *multiplicatively*, so that doubling their sizes leads to quadrupling their contribution to the block output. The simultaneous decrease of block multiplier by 4 then roughly keep their contribution invariant in size.

Large depth On the other hand, the optimal sublevel set has slope ≈ -1 when the depth is 2^{10} : Doubling the learning rate while halving the block multiplier has similar performance. This reflects the fact that ΔW^{l1} and ΔW^{l2} now interact *additively*.

Intermediate depths interpolate this phenomenon, as seen in the plot for depth 2^5 .

In the same heatmaps, one can see the optimal (learning rate, block multiplier) (in the $1/\sqrt{L}$ parametrization) shifts from the middle of the grid to the lower left as depth goes from 2^5 to 2^{10} , demonstrating the lack of hyperparameter transfer.

This change in slope is seen in relu networks as well, with or without layernorm.

Finally, we note that the $1/\sqrt{L}$ scaling still yields a $L \rightarrow \infty$ limit where the network still learns features as a whole, even though within each block this is no longer true. Thus, this is another reminder that mere "feature learning" does not imply "hyperparameter transfer"!

M.3 Classification of Parametrizations

These heatmaps already demonstrate that no parametrization of (global learning rate²¹, block multiplier) can transfer hyperparameters robustly, because any such parametrization can only *shift* the heatmaps but not *stretch* them, so one cannot "transfer" a sublevel set of one slope into a sublevel set of another slope.

But even if we allow learning rate to vary between layers in a block, no stable, faithful, nontrivial parametrization can avoid the linearization problem described above.

For simplicity, fix a positive-homogeneous nonlinearity and block depth 2^{22} . We consider the space of hyperparameters consisting of the learning rate for each of the layers in a block, as well as the block multiplier (one for each block); WLOG all weights are initialized $\Theta(1)$.²³ This yields a space of dimension $\text{blockdepth} + 1 = 3$.

Indeed, for this to happen, the weight update ΔW^{li} must be at least of order $\Omega(1)$ (size of initialization) for some i . But this would contribute a drift term to the block output $g^l = g^l(x^{l-1}; W^\bullet)$ that is as large as the noise term. This then implies that either the parametrization is unstable (if the block multiplier $L^{-\alpha}$ is $\Omega(1/L)$) or lacks feature diversity (if the block multiplier $L^{-\alpha}$ is $O(1/L)$).

For example, in a linear model,

$$L^\alpha \langle g^l \rangle = \langle W^{l2} W^{l1} x^{l-1} \rangle = \langle W_0^{l2} W^{l1} x^{l-1} \widehat{} \rangle + \langle W_0^{l2} W^{l1} x^{l-1} \rangle + \langle \Delta W^{l2} W^{l1} x^{l-1} \rangle.$$

$\langle W_0^{l2} W^{l1} x^{l-1} \widehat{} \rangle$ is independent and zero-mean across l (the noise term), while $\langle W_0^{l2} W^{l1} x^{l-1} \rangle + \langle \Delta W^{l2} W^{l1} x^{l-1} \rangle$ is correlated across l (the drift term). $\langle W_0^{l2} W^{l1} x^{l-1} \widehat{} \rangle$ is always $\Theta(1)$ because the W_0^{l2}, W_0^{l1} are. If ΔW^{l2} is $\Omega(1)$, then $\langle \Delta W^{l2} W^{l1} x^{l-1} \rangle = \Omega(1)$ as well, making the drift term as large as the noise term. If ΔW^{l1} is $\Omega(1)$, then $\langle W_0^{l2} \Delta W^{l1} x^{l-1} \rangle = \Omega(1)$, causing $\langle W_0^{l2} W^{l1} x^{l-1} \rangle = \langle W_0^{l2} W_0^{l1} x^{l-1} \rangle + \langle W_0^{l2} \Delta W^{l1} x^{l-1} \rangle$ to be $\Omega(1)$.²⁴

The same argument can be straightforwardly adapted to nonlinear MLPs (with mean subtraction) and arbitrary block depth ≥ 2 , and as well to general nonlinearities that are not necessarily positive-homogeneous, with hyperparameter space enlarged to include initialization.

M.4 So What is the Optimal Parametrization?

All of the above considerations suggest that *we are missing crucial hyperparameters in our consideration* when increasing the complexity of each block. Our study right now is akin to the naive study of the 1-dimensional hyperparameter space of the global learning rate in SP. Discovering these missing hyperparameters will be an important question for future work.

N Experiments

N.1 Verifying the Theory in the Linear Case

In Appendix C, we showed that a complete description of the training dynamics of linear networks can be formulated in terms of Γ and C . In this section, we provide empirical results supporting our theoretical findings. We first verify the finite-depth recursive formula for Γ in Lemma C.3 is the correct limit when the width goes to infinity, then proceed to show that the infinite-depth limit is the correct one.

²¹meaning, the learning tied across all layers in a block

²²but our arguments generalize trivially to arbitrary block depth ≥ 2

²³This is WLOG because the nonlinearities are homogeneous

²⁴One can also observe that if $\Delta W^{l1} = \Omega(1)$, then by symmetry the backward pass suffers the same problem. But for general block depth, this argument does not say anything about the middle layers, while the argument presented above implies that ΔW^{li} cannot be $\Omega(1)$ for any i .

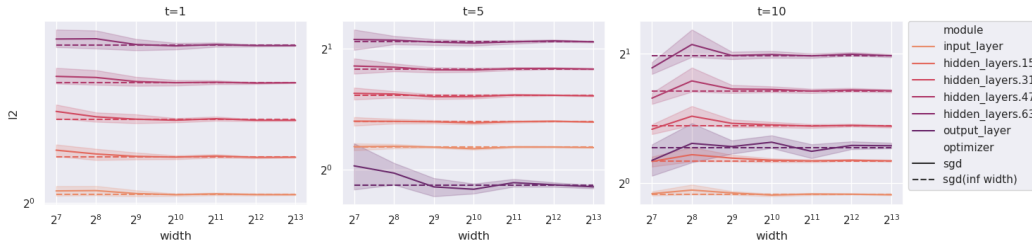


Figure 9: Trained linear network converges to its infinite width limit which is computed recursively based on Γ and C . Depth is fixed at 64, width varies between $2^7, 2^8, \dots, 2^{13}$. Networks are trained with SGD for 10 steps. The root mean square statistics (y -axis) at 1st, 5th and 10th steps are plotted using solid lines where the x -axis is the width. The root mean square values are computed on the outputs of some of the layers (including the input layer, output layer, and hidden layers at each quarter). The corresponding value for the infinite width is indicated with dashed lines.

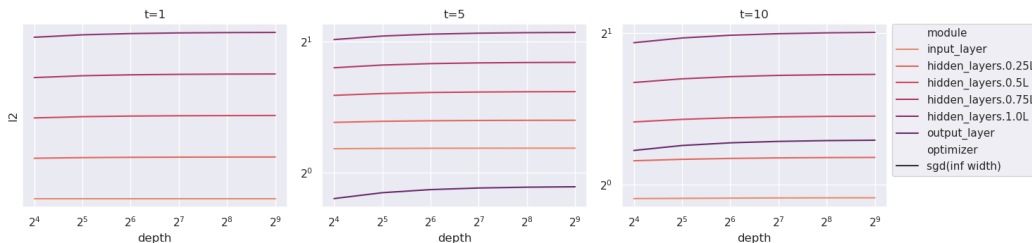


Figure 10: Under Depth- μP , infinite wide linear network training converges when increasing the depth. Infinite wide linear networks of depth $2^4, 2^5, \dots, 2^9$ are computed recursively based on Γ and C . The root mean square statistics (y -axis) at 1st, 5th and 10th steps are plotted across the depth (x -axis).

Infinite-width limit. In Figure 9, we train a series of 64-layer linear networks of width $2^7, 2^8, \dots, 2^{13}$ with 1, 5, 10 steps on MNIST, and plot the root mean square²⁵ of the layer outputs using solid lines. We also compute the infinite width limit of the corresponding statistics using the recursive formula for Γ and plot them as dashed horizontal lines. For clarity of the figure, we only plot the statistics of the input layer, output layer, and hidden layers of index 16, 32, 48, and 64. It is clear that as the width grows, the solid lines converge to the dashed lines consistently across the training steps. It indicates that our computation of the infinite width limit is correct.

Infinite-depth limit. We verify that the infinite *width* limit above converges when the *depth* grows. We consider linear networks of the same architecture but vary the depth from 2^4 to 2^9 . We again compute the root mean square values of the layer outputs using the recursive formula for Γ , and plot them in Figure 10 with depth being x -axis. For clarity of the figure, we only plot the statistics of the input layer, output layer, and hidden layers of index $L/4, L/2, 3L/4$, and L . One can observe that the statistics of the layer outputs converge quickly when the depth grows from 2^4 to 2^9 , which verifies our convergence result.

N.2 Hyperparameter Transfer

In this section, we provide empirical evidence to show the optimality of Depth- μP scaling and the transferability of some quantities across depth. We train vanilla residual network with block depth 1 (1 MLP layer in each residual block) on CIFAR-10 dataset using Adam optimizer, batch size 64, for

²⁵The root mean square of a vector $x = (x_1, \dots, x_n)$ is $\sqrt{\frac{\sum_{i=1}^n x_i^2}{n}}$, which is denoted as “l2” in Figures 9 and 10.

50 epochs (input and output layers are fixed). The network is parameterized as follows

$$x^l = x^{l-1} + a \times L^{-\alpha} \text{MS}(\phi(W^l x^{l-1})),$$

and the weights are trained with the rule

$$W^l \leftarrow W^l - \eta \times n^{-1} L^{-\gamma} Q_t^l(nL^\delta g_0, \dots, nL^\delta g_t),$$

where the learning rate η and the block multiplier a are the *hyperparameters*.²⁶ The values of α, γ depend on the parametrization of choice. For Depth- μ P, we have $\alpha = \gamma = 1/2$, and for standard parametrization, we have $\alpha = 0, \gamma = 1$.²⁷ In our experiments, we assume base depth 8, meaning that we replace L by $L/8$ in the parametrization above.

Learning rate transfer (η). In Figure 11, we show the training loss versus learning rate for depths 2^k , for $k \in \{3, 4, \dots, 10\}$. For Depth- μ P, a convergence pattern can be observed for the optimal learning rate as depth grows. Optimal learning rates for small depths (e.g. $L = 2^3$) exhibit a mild shift which should be expected, as our theory shows convergence in the large depth limit. However, starting from depth $L = 2^6$, the optimal learning rate is concentrated around 10^{-3} . With standard parametrization with learning rate scaling ($\alpha = 0, \gamma = 1$), the optimal learning rate exhibits a significant shift with depth and training loss degrades when the depth is too large. Here we have already set $a = 2^{-3}$, making the curves look better compared to the standard practice with $a = 1$. For standard parametrization without any depth scaling ($\alpha = \gamma = 0$), the optimal learning rate exhibits a significant shift as depth grows, suggesting that standard parametrization is not suitable for depth scaling. Additional figures with multiple time slices are provided in Appendix O.

Is feature learning sufficient for HP transfer? In Appendix F, we explained when and why hyperparameter transfer occurs. Precisely, to obtain HP transfer, one needs to classify all feature learning limits and choose the optimal one. We introduced the notion of feature diversity and showed that Depth- μ P is optimal in the sense that it maximizes feature diversity. To show that optimality is needed for HP transfer, we train a resnet with $(\alpha, \gamma) = (1, 0)$ which is also a feature learning limit. Figure 12 shows that in this case the learning rate exhibits a significant shift with depth. Interestingly, the constant η in this case seems to increase with depth, suggesting that the network is trying to break from the *ODE* limit, which is sub-optimal. Note that in Figure 11, with Depth- μ P we obtain better training loss compared to the *ODE* parametrization in Figure 12.

Do we still have transfer with LayerNorm (LN)? Our theory considers only Mean Substraction (MS), and Figure 11 shows the results with MS. To see whether LN affects HP transfer, we train resnets with the same setup as Figure 11 with absolute value non-linearity and LN applied to x^{l-1} before matrix multiplication with W^l (preLN). We keep MS after non-linearity although it can be removed since LN is applied in the next layer. Our results, reported in Figure 13 suggest that Depth- μ P guarantees learning rate transfer with LN as well.

Block multiplier transfer (a). In Figure 14, we investigate the stability of the hyperparameter a in Depth- μ P as depth increases. The results suggest that the optimal value of this constant converges as depth grows, which suggest transferability. Additional experiments with multiple time slices are provided in Appendix O.

N.3 What Happens in a Transformer?

Because transformers have block depth 2, as discussed in appendix M, we have plenty of reasons to suspect that no parametrization of (learning rate, block multiplier) will be able to robustly transfer hyperparameters across depth for transformers.

Here we do a large scale experiment using Megatron trained on Common Crawl and catalogue our observations.²⁸ In summary, in our particular setup (which should be close to most large

²⁶Note that η here is the constant, and the effective learning rate is given by $\eta n^{-1} L^{-\gamma}$.

²⁷In standard parametrization, there is generally no rule to scale the learning rate with depth, and the optimal learning rate is typically found by grid search. Here, we assume that in standard parametrization, the learning rate is scaled by L^{-1} to preserve faithfulness.

²⁸We train the models for 3900 steps, using cosine decay schedule with 500 warmup steps. We use a sequence length of 4096, batch size 256, resulting in approximately 4B tokens per training run.

language model pretraining), we see that the $1/\sqrt{L}$ scaling seems to transfer hyperparameters at the end of training (fig. 17(Right)). However, we also see that 1) deeper does worse in initial training (fig. 16(Left)), and 2) optimal hyperparameters scale like $\Theta(1)$ in the middle of training (fig. 17(Left)). Combined with the theoretical insights of appendix M, this leads us to conclude that while the $1/\sqrt{L}$ scaling can potentially be practically useful in transformer training, it is likely to be brittle to architectural and algorithmic changes, or even simple things like training time.

In fact, we observe that transformers are insensitive to the block multiplier a (fig. 15), so that the only relevant hyperparameter is really just learning rate. Thus, empirically measuring the scaling trend of the optimal learning rate, as done in modern large scale pretraining, can be a practically more robust way to transfer hyperparameters.

Here L is the number of transformer layers, each of which consists of an attention layer and an MLP layer (each of which has depth 2).

N.4 Feature Diversity

In this section, we empirically verify our claims about feature diversity exponent (Claims J.4 and J.5). We use the same setup as in the last section, i.e., we train deep residual networks of width $n = 256$ on CIFAR-10 dataset with Adam and batch size 64. In Figure 18, we compare two parametrizations, Depth- μ P ($\alpha = \gamma = 1/2$) and the ODE parametrization $(\alpha, \gamma) = (1, 0)$. We measure $\left\| \mathbf{x}_t^{\lfloor (\lambda + \epsilon)L \rfloor} - \mathbf{x}_t^{\lfloor \lambda L \rfloor} \right\| \stackrel{\text{def}}{=} d(\epsilon)$ at $t = 1000$ for the two parametrizations and varying depth. For each parametrization and depth L , we rescale function d by multiplying a constant c such that $c \cdot d(1/256) = 1$, and then plot the rescaled function $c \cdot d$ for a clean presentation. One can observe clearly that Depth- μ P has feature diversity exponent (almost) $1/2$ for any L , while the curves for ODE parametrization move from $\epsilon^{1/2}$ to ϵ when L grows. This exactly fits our theory that Depth- μ P maximizes the feature diversity, while other parametrizations (even with feature learning) have smaller feature diversity exponents that should go to 0 in the infinite depth limit.

Growth along with L and t . In Figure 19, we measure $d(\epsilon)$ at $t = 100, 500, 1000$, and rescale it by dividing additional $\epsilon^{0.5}$ and a constant c such that $\frac{d(1/256)}{c \cdot \epsilon^{0.5}} = 1$, and then plot the rescaled function $d/(c \cdot \epsilon^{0.5})$ for a clean comparison between d and $\epsilon^{0.5}$. We observe that for both Depth- μ P and ODE parametrization, the slopes of the curves grow along with L and t . The growth along t can be explained by the cumulative correlation between layers. The growth along L for ODE parametrization is because the independent components between nearby layers decrease when L grows. We do not have a clear understanding for the growth along L for Depth- μ P and we leave it as a future work.

Absolute value activation increases feature diversity. In Figure 20, we plot the same curves as in Figure 19 but comparing ReLU activation and absolute value activation under Depth- μ P. We observe that the slope of the curves for absolute value activation is smaller than ReLU activation. It matches our theory that absolute value activation increases feature diversity.

O Additional Experiments

O.1 Failure of Standard Parametrization at Large Depths

O.2 Experiments with Block Depth 2

Currently, our theory covers resnets with block depth 1, and our experiments confirm the theoretical findings. We conducted similar experiments for block depth 2 (i.e. the residual block consists of 2 fully connected layers) to see whether the learning rate transfers with Depth- μ P. The results are reported in Figure 22. The results show a significant shift in the learning rate which might indicate that as block depth increases, adjustments are needed to stabilize hyperparameters with depth.

O.3 Other experiments

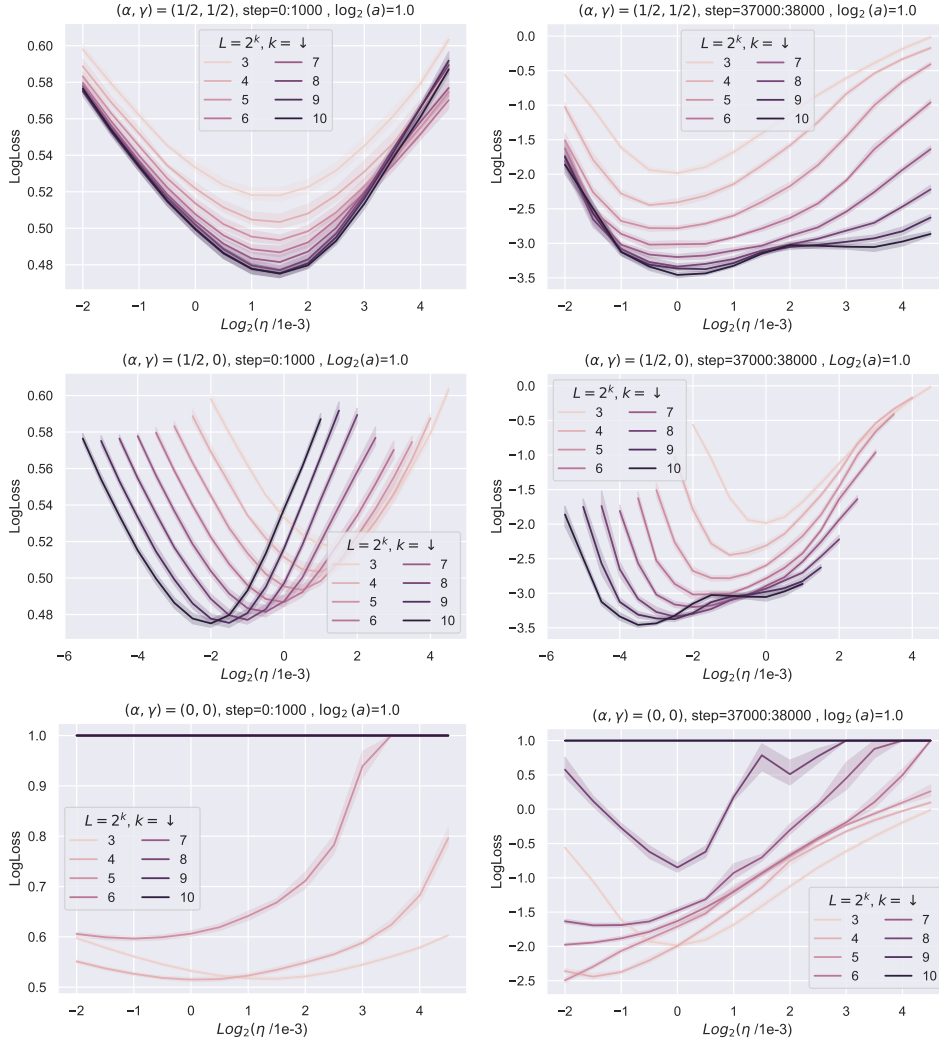


Figure 11: Train logloss versus learning rate for width $n = 256$ and varying depths. The network consists of MLP blocks (with block depth 1), trained for 50 epochs on CIFAR10 dataset using Adam. The batch size is fixed to 64. We tune the depth 2^3 network to obtain the optimal $(\log_2(a), \log_2(\eta/1e-3)) = (1, 0)$, and scale all deeper networks using 2^3 as base depth. The reader can check that the $L = 2^3$ curves in each columns are the same. We show the logloss versus the learning rate of the hidden layers (input/output layers fixed) for three parametrizations: Depth- μ P (**Top**), Scaling only the blocks (no LR scaling), i.e. $\gamma = 0$ (**Middle**), and Standard Parametrization without any scaling ($\alpha = \gamma = 0$) (**Bottom**). Each curve represents the average training loss over a time slice of 1000 steps for depths 2^k for $k \in \{1, 2, \dots, 10\}$. Confidence intervals are based on 5 seeds. The results show that Depth- μ P preserves the optimal learning rate while consistently improving the training loss as depth increases. If we only scale the blocks without scaling the LR ($\alpha = 1/2, \gamma = 0$) when training with Adam, the optimal learning rate shifts significantly with depth. With standard parametrization without any depth scaling (common practice), the results show a significant shift in the optimal learning rate as well. For SP, we cap the log loss at 1, which is why for depth $2^9, 2^{10}$, we have a black horizontal line at $\text{LogLoss} = 1$.

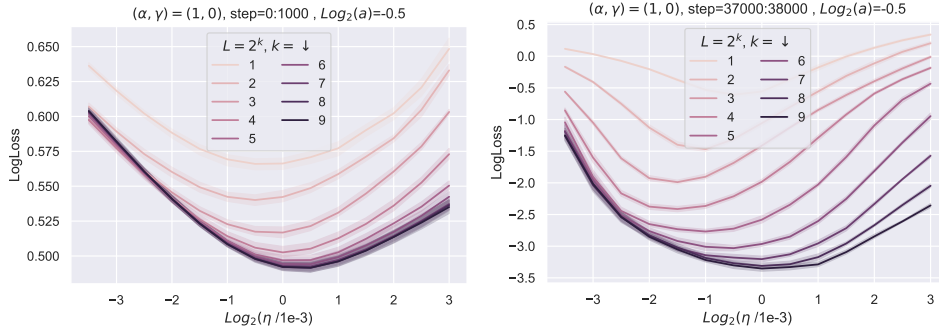


Figure 12: Same setup as fig. 11 for the parametrization $(\alpha, \gamma) = (1, 0)$ (the ODE limit).

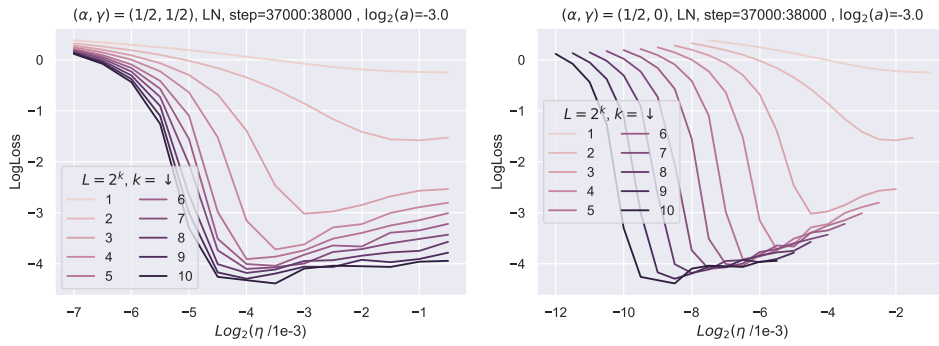


Figure 13: Same setup as Figure 11 with Abs non-linearity instead of ReLU and LayerNorm applied to x^{l-1} before matrix multiplication with W^l . We show the logloss versus the learning rate of the hidden layers (input/output layers fixed) for two parametrizations: Depth- μ P (**Left**) and scaling only the blocks without LR scaling ($(\alpha, \gamma) = (1/2, 0)$) (**Right**). The results show that Depth- μ P preserves the optimal learning rate while consistently improving the training loss as depth increases. If we only scale the blocks without scaling the LR ($\alpha = 1/2, \gamma = 0$) when training with Adam, the optimal learning rate shifts significantly with depth.

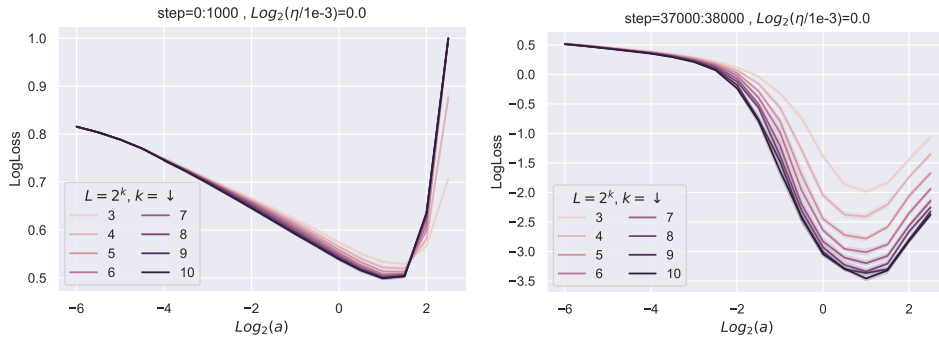


Figure 14: Train logloss versus block multiplier a for varying depths. Same training setup as in fig. 11. The results suggest that Depth- μ P stabilizes the hyperparameter a as depth increases.

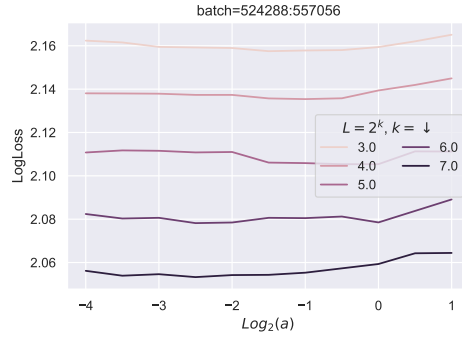


Figure 15: Modern transformers are insensitive to block multiplier a .

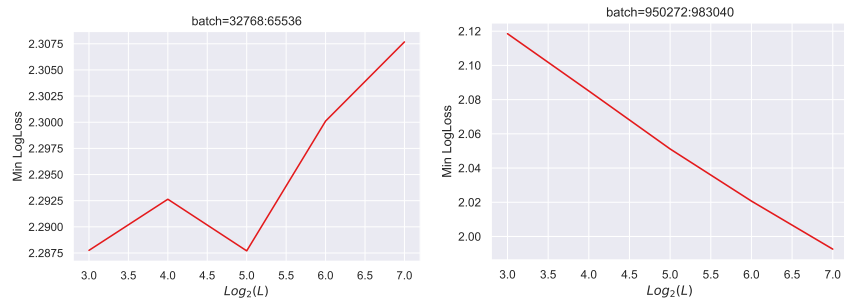


Figure 16: In (Megatron) Transformer trained on Common Crawl, deeper does worse initially (Left) but eventually does better (Right).

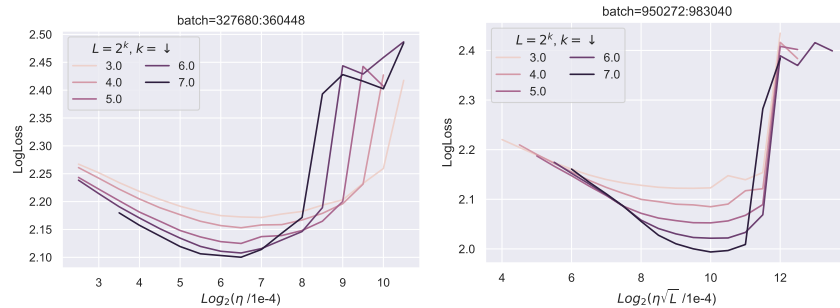


Figure 17: In the middle of (Megatron) transformer training, optimal learning rate is approximately invariant (Left), while at the end of training, it approximately scales like $1/\sqrt{L}$. However, the $1/\sqrt{L}$ scaling transfers the maximum viable learning rate better in either case.

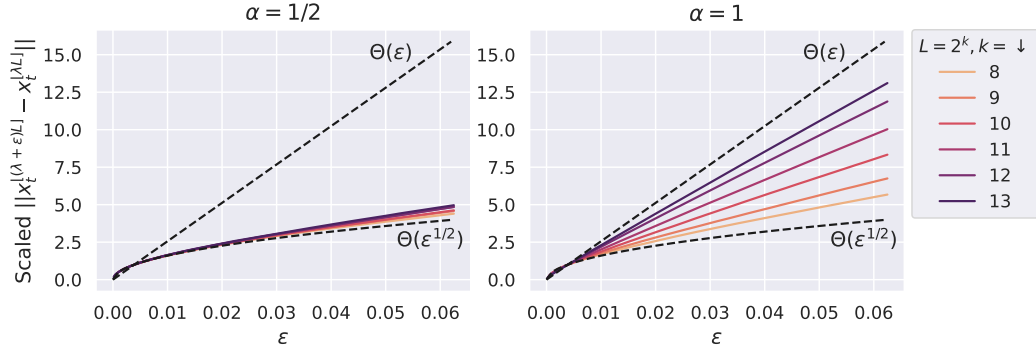


Figure 18: Difference between feature at layer $\lfloor \lambda L \rfloor$ and feature at layer $\lfloor (\lambda + \epsilon)L \rfloor$ as a curve of ϵ for width $n = 256$ and varying depths. For a clean presentation, each curve is scaled so it always passes $(1/256, 1)$. The feature diversity exponent κ depends on the growth of the curve when $L \rightarrow \infty$. For Depth- μ P (left), the curve is always close to $\epsilon^{1/2}$, meaning $\kappa = 1/2$. For ODE parametrization (right), the curve shifts from $\epsilon^{1/2}$ to ϵ when L grows, indicating its κ goes to 0 in the infinite depth limit.

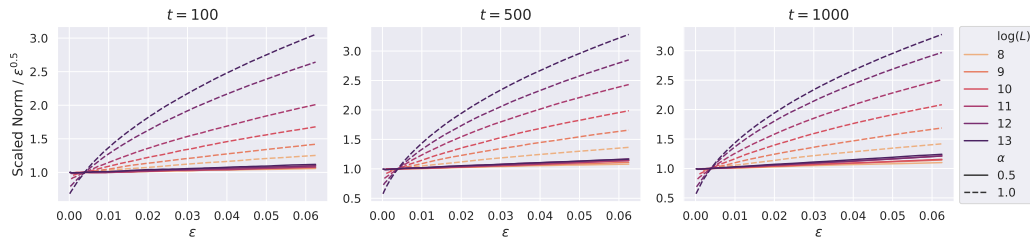


Figure 19: Same setup as Figure 18 but at step $t = 100, 500, 1000$, and each curve is scaled by dividing a constant and *additional* $\epsilon^{1/2}$ so it always passes $(1/256, 1)$. The curve indicating feature diversity exponent κ exactly $1/2$ should be a horizontal line at 1. For Depth- μ P ($\alpha = 0.5$), the curves are almost horizontal. For ODE parametrization ($\alpha = 1$), slopes of the curves are larger with larger L and larger t .

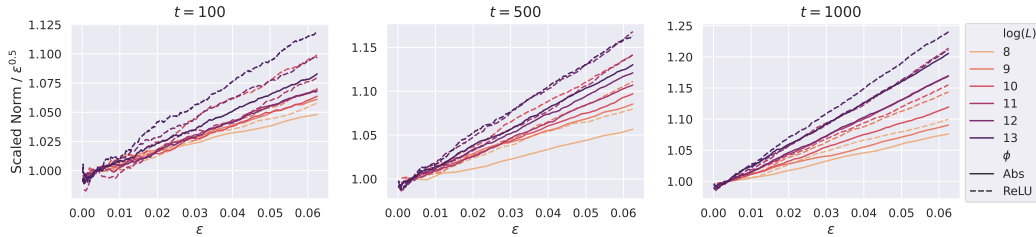


Figure 20: Same setup as Figure 19, but comparing Depth- μ P with ReLU activation and absolute value activation. Each curve is scaled by dividing a constant and $\epsilon^{1/2}$ so it always passes $(1/256, 1)$. The curve indicating feature diversity exponent κ exactly $1/2$ should be a horizontal line at 1. For both activations, slopes of curves are small, but growing along with L and t . The slopes with absolute value activation ($\phi = \text{Abs}$) are slower than the slopes with ReLU activation ($\phi = \text{ReLU}$), indicating feature diversity is higher with absolute value activation.

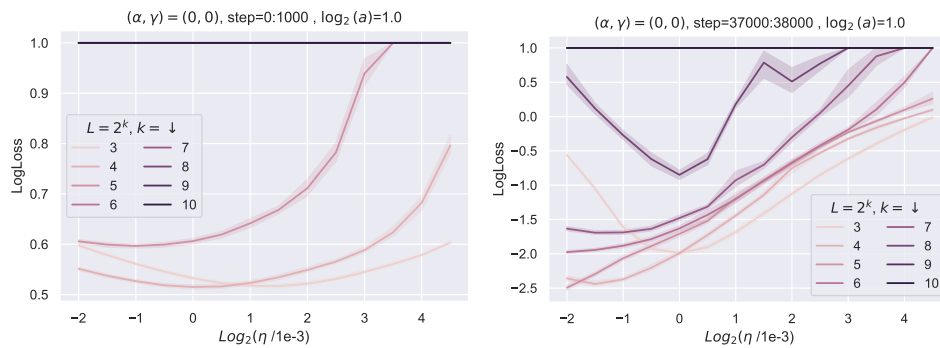


Figure 21: Training with Standard Parametrization fails at large depth due to numerical issues.

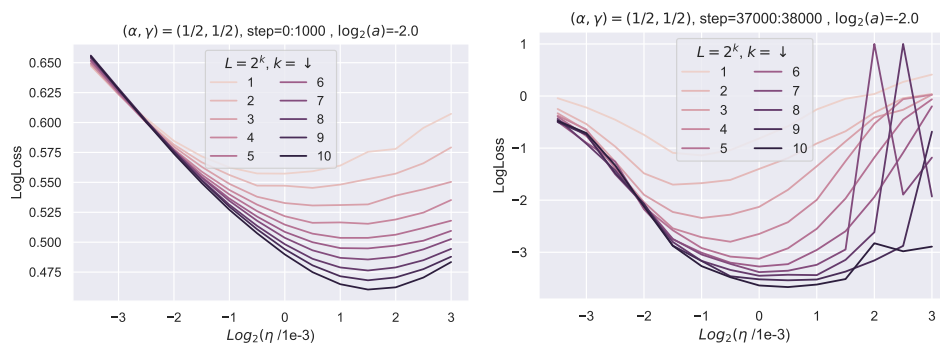


Figure 22: Same setup as Figure 11, with block depth 2 instead.

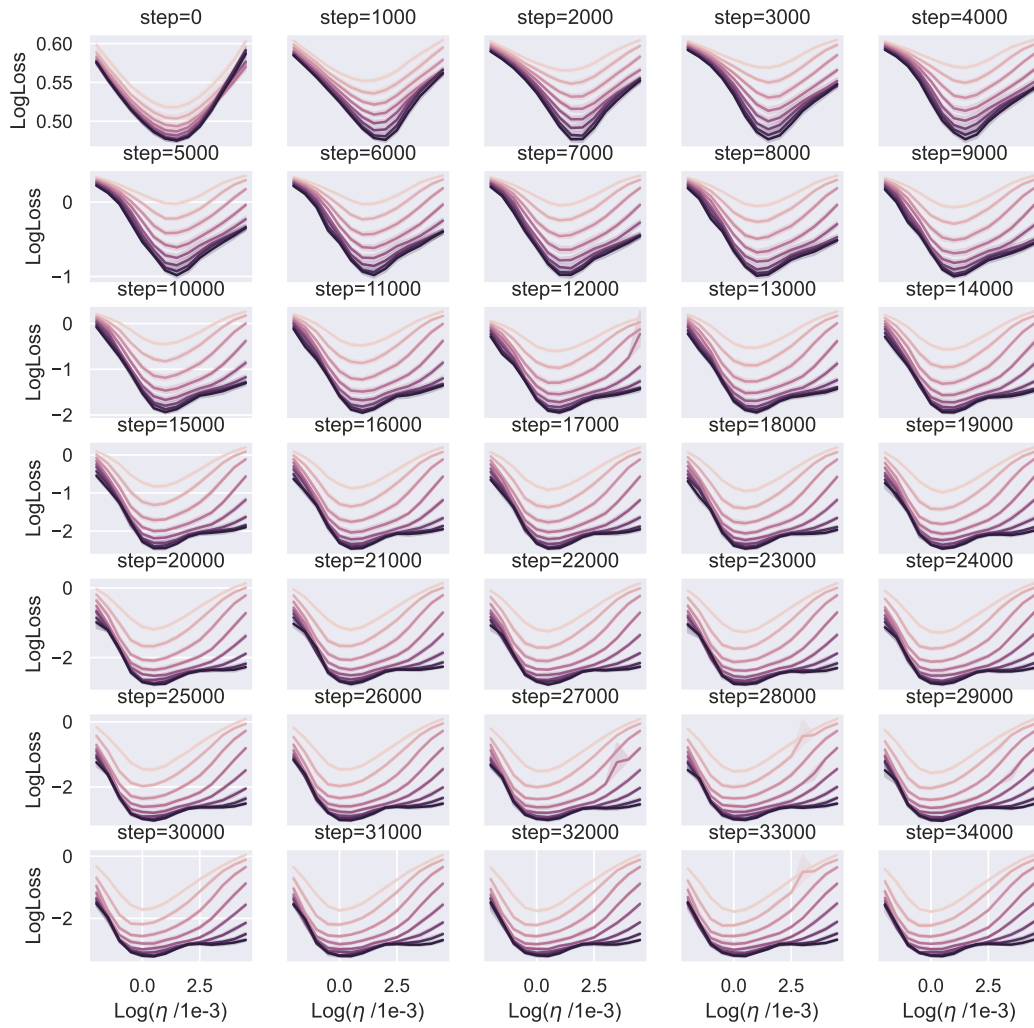


Figure 23: Same as fig. 11 (**Up**, Depth- μ P) with multiple time slices.

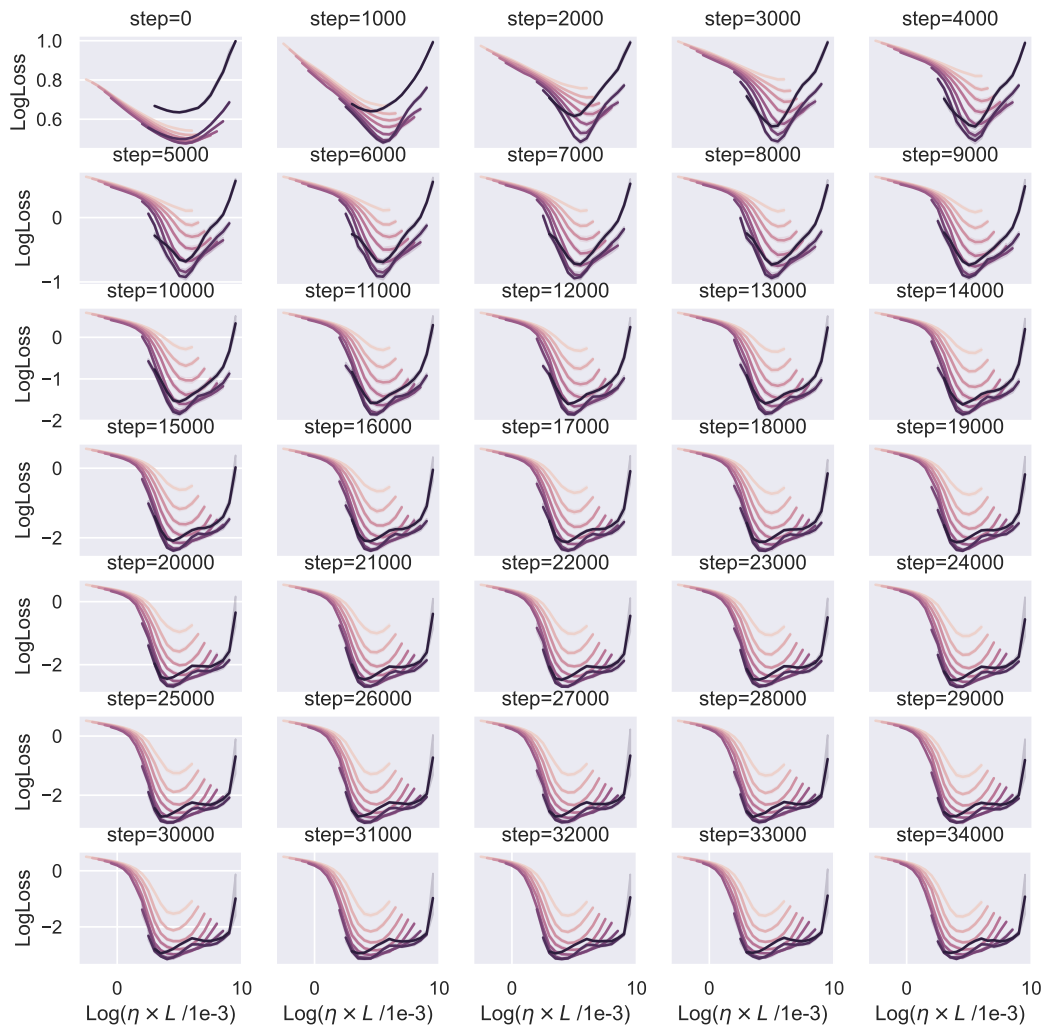


Figure 24: Same as fig. 11 (**Middle**, Standard Parametrization with $\gamma = 1$) with multiple time slices.

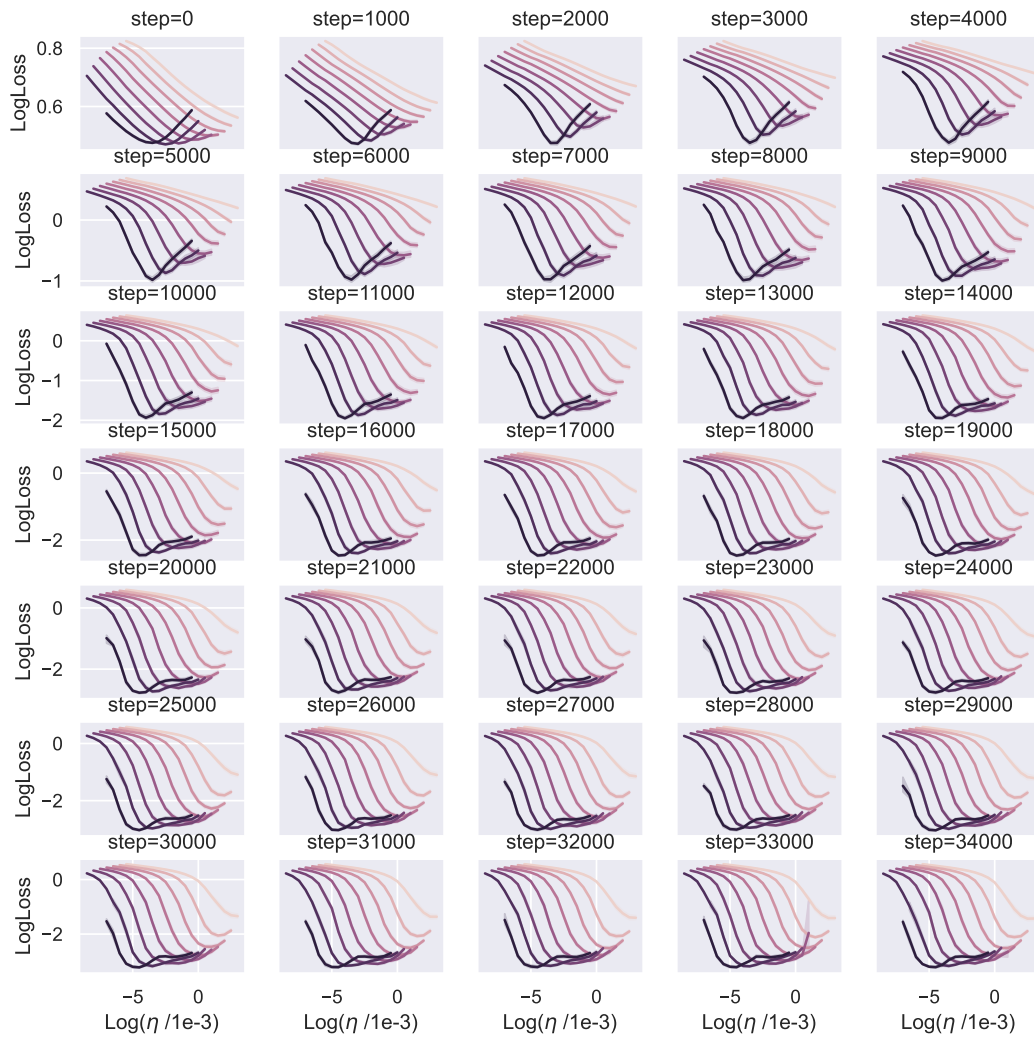


Figure 25: Same as fig. 11 (**Bottom**, Standard Parametrization with no scaling, $\alpha = 0, \gamma = 0$) with multiple time slices.

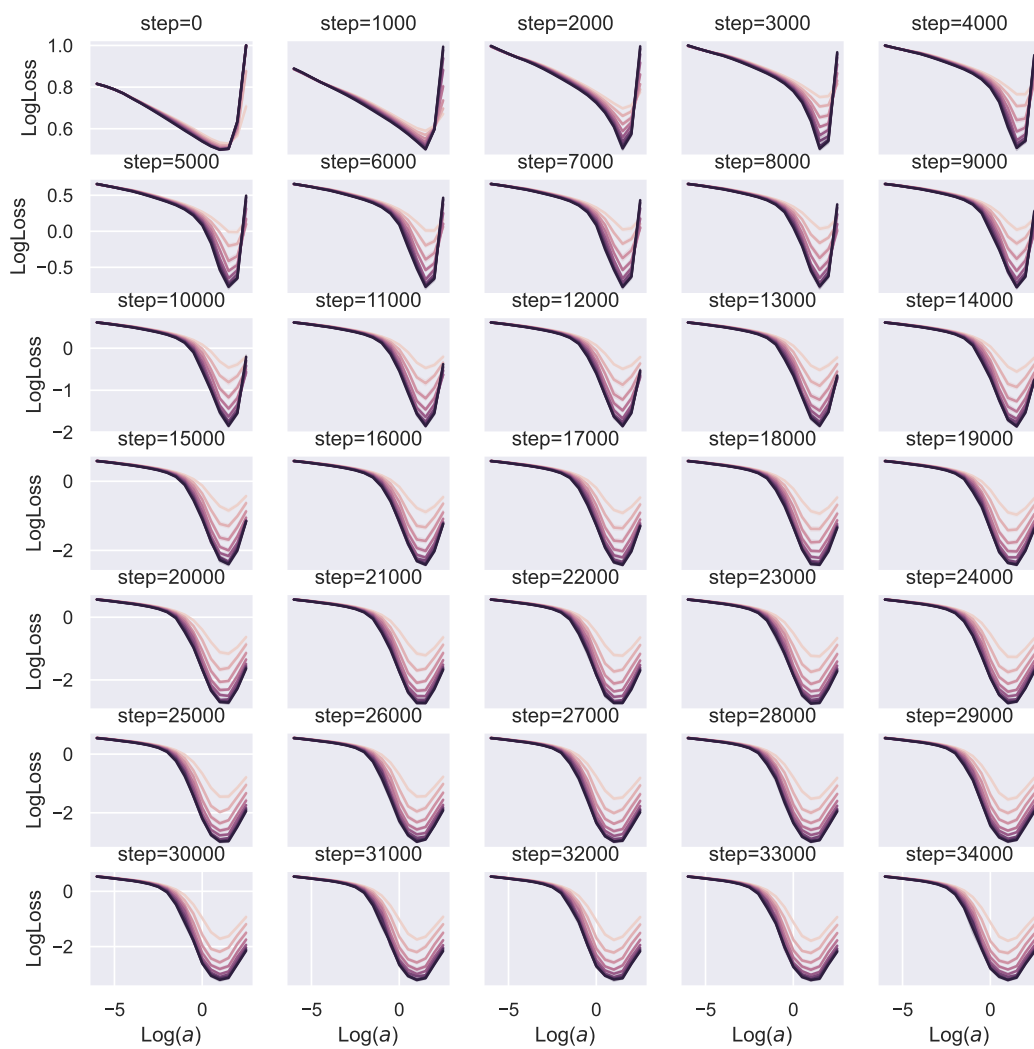


Figure 26: Same as fig. 14 with multiple time slices.