

# UNDERSTANDING THE ROLE OF DEPTH IN THE NEURAL TANGENT KERNEL FOR OVERPARAMETERIZED NEURAL NETWORKS

**William St-Arnaud**

Université de Montréal & Mila [william.st-arnaud@umontreal.ca](mailto:william.st-arnaud@umontreal.ca)

**Margarida Carvalho**

Université de Montréal & Mila [carvalho@iro.umontreal.ca](mailto:carvalho@iro.umontreal.ca)

**Golnoosh Farnadi**

McGill University & Mila  
[farnadig@mila.quebec](mailto:farnadig@mila.quebec)

## ABSTRACT

Overparameterized fully-connected neural networks have been shown to behave like kernel models when trained with gradient descent, under mild conditions on the width, the learning rate, and the parameter initialization. In the limit of infinitely large widths and small learning rate, the kernel that is obtained allows to represent the output of the learned model with a closed-form solution. This closed-form solution hinges on the invertibility of the limiting kernel, a property that often holds on real-world datasets. In this work, we analyze the sensitivity of large ReLU networks to increasing depths by characterizing the corresponding limiting kernel. Our theoretical results demonstrate that the normalized limiting kernel approaches the matrix of ones. In contrast, they show the corresponding closed-form solution approaches a fixed limit on the sphere. We empirically evaluate the order of magnitude in network depth required to observe this convergent behavior, and we describe the essential properties that enable the generalization of our results to other kernels.

## 1 INTRODUCTION

Machine learning approaches have demonstrated a remarkable ability in helping solve a multitude of problems across a large span of tasks. Whether the task is classification, prediction, or image generation, some variants of the combination of a neural network plus gradient descent have managed to achieve superhuman ability in some instances (LeCun et al., 1998; Silver et al., 2018; Vaswani et al., 2017). In recent years, a particular observation has been made regarding neural networks that are overparameterized. While previous common thought indicated that these large networks would fall prey to overfitting, this conclusion is being challenged empirically (Belkin, 2021). This phenomenon is not quite well understood and it has led to works analyzing the learning dynamics of overparameterized models updated with gradient descent (Liu et al., 2020; 2022; Jacot et al., 2018). While these works provide insight as to the learning dynamics of fully-connected neural networks that are overparameterized, Jacot et al. (2018), in particular, offer a closed-form solution to the gradient flow based on a kernel that is recursively computed. This result offers the possibility of approximately predicting the output of an overparameterized neural network learned through gradient descent, without explicitly training the model. This comes at the cost of computing a kernel over a particular dataset, which involves the computation of expectations.

---

We focus on fully-connected ReLU networks, which allows us to both speed up the computation of the kernel and obtain a more interpretable closed-form solution. Leveraging this more interpretable closed-form solution, we study the role of depth in the limiting kernel of infinitely wide fully-connected ReLU networks. Our contributions address two central aspects of the effect of increasing depths: 1) the convergence of the kernel, established in Proposition 4, and 2) the limiting solution to the output of a fully-connected ReLU network under infinitely wide hidden layers and infinitely small learning rate via Theorem 3. Our results apply to arbitrary data with support on the sphere and, in contrast to previous literature, do not require any assumptions on the spectrum of the Hermite expansion or the Mercer decomposition of the kernel. In contrast to Hanin & Nica (2020), we study the deterministic limit of the neural tangent kernel when the width is much larger than the depth: while we allow the depth to increase to infinity, the rate at which it does so is much slower than the widths of the hidden layers of a neural network. We offer a summary of results in relation to ours in table 1 in Appendix D.

## 2 RELATED WORK

The learning dynamics of overparameterized neural networks have been extensively studied through the lenses of kernel methods and Hessian-based analysis. A key development in this area is the neural tangent kernel (NTK), introduced by Jacot et al. (2018), which describes how infinitely wide fully-connected neural networks trained with gradient descent evolve linearly in function space. The NTK framework formalizes how, under common assumptions—particularly Gaussian initialization and wide-layer limits—neural networks behave similarly to kernel methods during training. Arora et al. (2019b) extend the NTK to convolutional neural networks (CNN), further highlighting the framework’s capacity to capture learning dynamics.

Subsequent work has reinforced the kernel-based interpretation of training dynamics through analyses of the Hessian matrix. Liu et al. (2020; 2022) and Belkin (2021) demonstrate that the loss landscape of overparameterized neural networks often exhibits near-linearity, with low-curvature regions and small-norm Hessians, supporting the NTK-based approximation. These findings suggest that network outputs are relatively stable during training, especially for wide architectures with standard initialization. Lee et al. (2020) provide an in-depth empirical analysis of NTK models and their performance compared to finite-width neural networks of various architectures (e.g. fully-connected, CNNs). One key observation is that NTKs often outperform finite-width networks, yet are usually surpassed by conventional CNNs.

The NTK typically requires Monte Carlo estimation of expectations over Gaussian distributions, especially when nonlinearities from activations are involved. This reliance on sampling can be computationally expensive and introduces variance in the resulting kernel evaluations. However, closed-form expressions have been derived for certain activation functions, such as ReLU and leaky ReLU (Tsuchida et al., 2018). Additionally, while prior work has largely focused on width, less attention has been paid to how depth affects NTK sensitivity to initialization and the associated multiplicity in network outputs. Among the few works addressing this dimension, Bietti & Bach (2020) show that for the uniform measure on the sphere, the reproducing kernel (c.f. reproducing kernel Hilbert space or RKHS) leads to the same representation power regardless of the network depth. This raises the question of the significance of depth for the NTK. Further insights are provided by Nguyen et al. (2021), who derive an asymptotic lower bound on the smallest eigenvalue of the NTK through the Hermite expansion of the kernel. This result can be used to derive bounds on the generalization of the model (Arora et al., 2019a) and gives a better grasp on understanding the role that depth plays in both convergence and generalization. Murray et al. (2023) characterize the full spectrum of the NTK via the Hermite expansion for arbitrary datasets on the sphere. They recover an empirical observation that the eigenvalues of the NTK follow power law decay with respect to the size of the training set; Li et al. (2024) extend this line of results to general domains beyond the sphere. In addition to the previous result from Jacot et al. (2018) regarding the convergence to kernel regression in the infinite width case, they also establish a uniform convergence bound to the NTK regressor for the output of the trained

model. This improves upon the pointwise convergence bounds found in Lee et al. (2019); Arora et al. (2019b); Allen-Zhu et al. (2019). As in the initial paper by Jacot et al. (2018), the limiting kernel is observed to be deterministic. This is contrasted with the work of Hanin & Nica (2020), where the authors characterize the NTK directly through the ratio of its second and first moment. They obtain a limiting result showing that if the ratio of depth to width is unbounded, the NTK has a much higher variance than mean, implying that it is highly stochastic.

### 3 NOTATION

To highlight the role of the depth  $L$  of a neural network, we denote elements that depend on layer  $l \in \{1, \dots, L\}$  with a superscript  $(l)$ , e.g.  $\kappa^{(l)}$ . To emphasize the special role of the last layer in the analysis, the superscript  $(L)$  is sometimes used in mathematical objects to remind the reader of its link to a neural network of depth  $L$ , e.g.  $\kappa^{(L)}$  refers to a network of depth  $L$  while it can also stand on its own as a mathematical object. The width of a layer  $l$  is denoted by  $n_l$ , with  $n_0$  referring to the input dimension. Dependence on a time parameter  $t$  and the size  $n$  of a dataset is sometimes indicated with a subscript to emphasize their importance, but is otherwise omitted. A dataset of size  $n$  is denoted  $X$ , understood as an  $n \times n_0$  matrix, where each row  $i$  is written as  $x_i^\top$  (i.e. the  $x_i$ 's are column vectors). Activation functions are denoted by  $\sigma$ , while uppercase  $\Sigma$  is reserved for computing "covariances" (see Definition 1). The limiting deterministic kernels of Jacot et al. (2018) are represented using  $\Theta_\infty^{(L)}$ , and  $\bar{\Theta}_\infty^{(L)}$  for their normalized version (see Definition 4); the notation  $\kappa$  is used when referring to general kernels and  $\bar{\kappa}$  for the normalized version. For the sake of simplicity, we consider neural networks with one-dimensional outputs (i.e.,  $n_L = 1$ ). Therefore, kernels refer in this context to functions  $\mathbb{R}^{n_0} \times \mathbb{R}^{n_0} \rightarrow \mathbb{R}_+$ . We also write  $\Theta(A)$  for the component-wise application of a kernel  $\Theta$  to the entries of a matrix  $A$ . Specifically,  $\kappa(XX^\top)$  denotes applying the kernel to all pairwise dot products in  $X$ , where  $XX^\top$  is the matrix containing those dot products. Similarly, for any function  $g: \mathbb{R} \rightarrow \mathbb{R}$ , the entry-wise application to a matrix  $A$  is denoted by  $g(A)$ . The notation  $A \leftarrow_{i,j} A'$  refers to the matrix obtained by replacing column  $i$  of matrix  $A$  with column  $j$  of matrix  $A'$ . The vector of ones of length  $n$  is denoted  $\mathbf{1}_n$ . Finally, the sphere of dimension  $n_0 - 1$  is denoted  $S^{n_0-1}$ .

### 4 BACKGROUND ON THE NTK AND OVERPARAMETERIZATION

Jacot et al. (2018) show that, under overparameterization and i.i.d. standard normal weight initialization, a fully-connected neural network of arbitrary depth  $L$  exhibits learning dynamics that converge to those of kernel gradient flow in the infinite-width limit. They also provide a recursive formula to compute the kernel  $\Theta_\infty^{(L)}$  to which gradient descent converges (see Theorems 1 and 2 from Jacot et al. (2018)). However, evaluating  $\Theta_\infty^{(L)}$  relies on computing high-dimensional expectations and can potentially also be subject to sample inefficiency in the approximation, motivating the search for a more readily computable kernel.

To make the kernel more practical, one may ask whether an efficient closed-form expression can be derived for particular activation functions. Of particular interest is the representation of  $\Theta_\infty^{(L)}$  in closed-form when using ReLU activations (Tsuchida et al., 2018), due their empirical popularity and methodological appeal in theoretical analysis. To this end, let us first introduce important definitions and recap the recursive formulation of  $\Theta_\infty^{(L)}$  by Jacot et al. (2018).

**Definition 1** ((mean) Covariance of neurons  $\Sigma^{(l)}$ ). *Let  $x$  and  $x'$  be two inputs in  $\mathbb{R}^{n_0}$ . The covariances of neurons from inputs  $x$  and  $x'$  at each layer  $l$  are defined recursively as*

$$\Sigma^{(1)}(x, x') := \frac{1}{n_0} x^\top x', \quad \Sigma^{(l+1)}(x, x') := \mathbb{E}_{f \sim \mathcal{N}(0, \Sigma^{(l)})} [\sigma(f(x)) \sigma(f(x'))]$$

where  $f \sim \mathcal{N}(0, \Sigma^{(l)})$  is an infinite vector indexed through the notation  $f(x)$  and  $f(x')$  and each vector  $(f(x), f(x'))^\top \sim \mathcal{N}(0, \Sigma^{(l)}(x, x'))$ . We also define the variant of  $\Sigma^{(l)}$  where we replace  $\sigma$  with its derivative  $\dot{\sigma}$ :

$$\dot{\Sigma}^{(l+1)}(x, x') := \mathbb{E}_{f \sim \mathcal{N}(0, \Sigma^{(l)})} [\dot{\sigma}(f(x)) \dot{\sigma}(f(x'))].$$

---

**Definition 2** (Neural tangent kernel (NTK)). *For inputs  $x$  and  $x'$ , the neural tangent kernel of the neural network  $f(\cdot; \theta)$  with parameters  $\theta \in \mathbb{R}^P$  is given by*

$$\Theta^{(L)}(x, x') = \sum_{p=1}^P \frac{\partial f(x; \theta_p)}{\partial \theta_p} \otimes \frac{\partial f(x'; \theta_p)}{\partial \theta_p}.$$

**Theorem 1** (Jacot et al. (2018)). *Suppose we have a fully-connected neural network of depth  $L$  with non-linear activation. In the limit as layer widths  $n_1, \dots, n_{L-1} \rightarrow \infty$ , the neural tangent kernel (see Definition 2)  $\Theta^{(L)}$  converges in probability to a deterministic limiting kernel:*

$$\Theta^{(L)} \rightarrow \Theta_{\infty}^{(L)} \otimes I_{n_L},$$

where  $\Theta_{\infty}^{(l)}$  is defined recursively by

$$\begin{aligned} \Theta_{\infty}^{(1)}(x, x') &:= \Sigma^{(1)}(x, x') \\ \Theta_{\infty}^{(l+1)}(x, x') &:= \dot{\Sigma}^{(l+1)}(x, x') \Theta_{\infty}^{(l)}(x, x') + \Sigma^{(l+1)}(x, x'). \end{aligned}$$

We remark that, although we assume  $n_L = 1$  for the sake of simplicity, Theorem 1 is stated in its general form for any output dimension  $n_L \in \mathbb{N}$ . We also note that this is the version of the theorem without biases ( $\beta = 0$  in the context of Jacot et al. (2018)). This theorem is key in the convergence results obtained in the next section (Proposition 4 and Theorem 2). We are now ready to state the simplified formula for positively correlated inputs (Proposition 1).

**Proposition 1.** *For ReLU activation and perfectly positively correlated inputs  $x$  and  $x'$ , i.e.  $\rho = 1$ , it holds that*

$$\Sigma^{(L)}(x, x') = \frac{1}{n_0 2^{L-1}} \|x\|_2 \|x'\|_2, \quad \dot{\Sigma}^{(L)}(x, x') = \frac{1}{2}$$

and

$$\Theta_{\infty}^{(L+1)}(x, x') = \frac{1}{2} \Theta_{\infty}^{(L)}(x, x') + \frac{1}{n_0 2^L} \|x\|_2 \|x'\|_2 = \frac{L+1}{n_0 2^L} \|x\|_2 \|x'\|_2.$$

*Proof sketch.* Note that  $\frac{x^\top x'}{\|x\| \|x'\|} = 1$  and the product  $\sigma^2(f(z))$ , where  $z \sim \mathcal{N}(0, 1)$ , follows a squared rectified gaussian distribution, and that  $\mu = 0$  implies  $x^\top x' \geq 0$  with probability  $\frac{1}{2}$ .  $\square$

**Definition 3** (Correlation coefficient of  $\Sigma^{(L)}(x, x')$ ). *The correlations of neurons from inputs  $x$  and  $x'$  are defined as*

$$\rho^{(L)}(x, x') := \frac{\Sigma^{(L)}(x, x')}{\sqrt{\Sigma^{(L)}(x, x) \Sigma^{(L)}(x', x')}}.$$

Note that  $\rho^{(L)}(x, x') \in [-1, 1]$ .

**Proposition 2** (Arora et al. (2019b)<sup>1</sup>). *For datapoints  $x$  and  $x'$  with  $\rho \in [-1, 1[$ , it holds that*

$$\begin{aligned} \rho^{(L+1)}(x, x') &= \frac{\sqrt{1 - (\rho^{(L)}(x, x'))^2}}{\pi} \\ &\quad + \frac{\rho^{(L)}(x, x') \arcsin \rho^{(L)}(x, x')}{\pi} \\ &\quad + \frac{1}{2} \rho^{(L)}(x, x') \\ \dot{\Sigma}^{(L+1)}(x, x') &= \frac{\arcsin \rho^{(L)}(x, x')}{2\pi} + \frac{1}{4} \end{aligned}$$

and

$$\Theta_{\infty}^{(L)}(x, x') = \|x\|_2 \|x'\|_2 \Theta_{\infty}^{(L)} \left( \frac{x}{\|x\|_2}, \frac{x'}{\|x'\|_2} \right)$$

for a fully-connected neural network with ReLU activation.

---

<sup>1</sup>See also Cho & Saul (2009) for the complete derivation.

With these results, we achieved our goal of obtaining a closed-form expression for the  $\Theta_\infty^{(L)}$  corresponding to an overparametrized (infinite-width), fully-connected ReLU network with no biases. This, in turn, allow us to aim at characterizing the output of such neural network, as done in the rest of the section. Indeed, from Proposition 2 and Proposition 2 from Jacot et al. (2018), we can immediately observe a few facts regarding the input data:

- case a) If all datapoints lie on the unit sphere  $S^{n_0-1}$ , the NTK is invertible for  $L \geq 2$  (Proposition 2 from Jacot et al. (2018)).
- case b) If all datapoints are pairwise not colinear, i.e.  $x_i^\top x_j < \|x_i\|_2 \|x_j\|_2$  for  $i \neq j$ , then the NTK is invertible (for  $L \geq 2$ ) since we can project them to different points on the sphere through the canonical projection.
- case c) If we project points from  $\mathbb{R}^{n_0}$  to the sphere  $S^{n_0-1}$  through stereographic projection, the embedding of the datapoints satisfies  $x_i^\top x_j = 1$  for all  $x_i, x_j$  in the dataset.

If one of these cases holds, the following proposition provides a closed-form expression for the approximation of the output of a fully-connected neural network.

**Proposition 3** (Jacot et al. (2018)). *Let  $X$  be a dataset of size  $n$  (with entries  $x_i^\top$ ) and let  $f^*$  and  $f_0$  respectively refer to the learned function and the neural network after the initialization. If the limiting kernel  $\kappa = \Theta_\infty^{(L)}(XX^\top)$  is invertible, the output of the neural network converges to*

$$f_\infty(x) = f_0(x) + \kappa_x^\top \kappa^{-1}(y^* - y_0),$$

where

$$\kappa_x = \Theta_\infty^{(L)}(xX^\top), \quad (y^*)_i = f^*(x_i), \quad (y_0)_i = f_0(x_i), \quad i = 1, \dots, n$$

as time  $t \rightarrow \infty$ , i.e. the number of gradient descent updates increases.

Motivated by Proposition 3 and the requirement to have an invertible  $\kappa$ , we identify **two** regimes of generalization: datapoints can lie on either a **1)** non-compact manifold (i.e.  $\mathbb{R}^{n_0}$ ) or **2)** a compact manifold (i.e.  $S^{n_0-1}$ ). The compact regime results in a simplifying assumption for the analysis that follows in this section. Note that one can project any dataset without any pair of colinear datapoints in  $\mathbb{R}^{n_0}$  on  $S^{n_0-1}$  using the canonical projection. The kernel  $\kappa$  will thus be invertible. If colinear points exist, a stereographic projection on  $S^{n_0-1}$  will result in an invertible  $\kappa$ .<sup>2</sup>

## 5 LIMITING KERNEL AS DEPTH INCREASES

While in the previous section, the depth  $L$  is fixed and the width goes to infinity, no mention is made of the effect of increasing both the depth and the width. Such insights into the additional effect of depth would provide a tangible frontier for the representation power of fully-connected neural networks and their generalization capabilities. In this section, we describe how the term  $\kappa_x$  from Proposition 3 approaches a fixed limit for each  $x$  as  $L \rightarrow \infty$ . The NTK will also approach this limit when  $L \rightarrow \infty$ , with  $L \in o(\min_{l=1, \dots, L-1} n_l)$ . Note that this setting is different from Hanin & Nica (2020), where the ratio of depth to width can be arbitrary; interestingly, when the depth grows faster than the width, there is no convergence to a deterministic limit for the NTK and it is stochastic. In this section, **we always assume ReLU activation** for  $\Theta_\infty^{(L)}$  and  $\Sigma^{(L)}$ . For a concise list that summarizes the assumptions made in this section, we refer the reader to Appendix A.

The following lemma demonstrates that  $\rho$  converges to 1 for each pair of datapoints as  $L$  goes to infinity. This result is a key ingredient in the propositions and theorems that follow.

**Lemma 1** (Convergence of  $\rho^{(L)}$ ). *If  $\rho^{(1)}(x, x') \in ]-1, 1[$ , then  $\rho^{(L)}(x, x') \rightarrow 1$  as  $L \rightarrow \infty$ .*

In the equation of Proposition 3, the terms  $\kappa_x$  and  $\kappa$  can both be normalized by a scalar and the resulting vector-matrix product is left unchanged. Specifically, if  $\Theta_\infty^{(L)}$  is normalized such that its diagonal elements are all equal to 1, some immediate results follow from Propositions 1 and 2. These results are shown in Proposition 4 and Theorem 2.

<sup>2</sup>In the context of learning the parameters of a neural network, we assume that one first projects onto the sphere and then fixes the projected data during the training phase.

**Definition 4** (Normalization of the  $\Theta_\infty^{(L)}$  kernel). *For  $x, x' \in S^{n_0-1}$ , the normalized version of  $\Theta_\infty^{(L)}$  is defined by*

$$\bar{\Theta}_\infty^{(L)}(x, x') = \frac{n_0 2^{L-1} \Theta_\infty^{(L)}(x, x')}{L}.$$

**Definition 5.** *We define the function  $h : [-1, 1] \rightarrow \mathbb{R}$  as*

$$h(z) = \frac{z \arcsin(z)}{\pi} + \frac{\sqrt{1-z^2}}{\pi} + \frac{z}{2}.$$

Using the definitions above, we state the Proposition 4 and the Theorem 2. The proofs can be found in Appendix B.

**Proposition 4** (Alternative formulation of  $\bar{\Theta}_\infty^{(L)}$ ). *The equality*

$$\bar{\Theta}_\infty^{(L+1)}(x, x') = \frac{L}{L+1} h'(\rho^{(L)}(x, x')) \bar{\Theta}_\infty^{(L)}(x, x') + \frac{1}{L+1} h(\rho^{(L)}(x, x'))$$

*holds  $\forall x, x' \in S^{n_0-1}$ . Moreover, the values in the normalized kernel are all found in the interval  $[0, 1]$ .*

**Theorem 2** (Convergence of  $\bar{\Theta}_\infty^{(L)}$ ). *For any  $x, x' \in S^{n_0-1}$ , the value  $\bar{\Theta}_\infty^{(L)}(x, x')$  strictly increases to 1 as  $L \rightarrow \infty$ .*

The result above can be taken to be a major obstacle to the analysis of  $\Theta_\infty^{(L)}(x^\top X^\top) \left( \Theta_\infty^{(L)}(XX^\top) \right)^{-1}$  since the positive determinant of  $\Theta_\infty^{(L)}$  converges to 0. However, Theorem 3, which is one of our key contributions, demonstrates that for a fixed  $x$ , the term  $\Theta_\infty^{(L)}(x^\top X^\top) \left( \Theta_\infty^{(L)}(XX^\top) \right)^{-1}$  converges to some limit as  $L$  increases to infinity. The proof of this theorem requires a function defined in Definition 6 and whose key properties are provided as Proposition 5. The required background on rough differential equations is provided in Appendix C.

**Definition 6.** *We define the function  $\psi_d$  for  $d \in \mathbb{R}^+$  as*

$$\psi_d(z) = \begin{cases} \frac{1}{1 + \exp\left(\frac{-2z}{d(1-z^2)}\right)} & \text{if } z \in ]-1, 1[ \\ 1 & \text{if } z = 1 \\ 0 & \text{if } z = -1. \end{cases}$$

**Proposition 5.** *The function  $\psi_d$  has the following key properties on  $[-1, 1]$ :*

$$\psi_d(-1) = 0 \tag{1}$$

$$\psi_d(1) = 1 \tag{2}$$

$$\psi_d \in \mathcal{C}^\infty \tag{3}$$

$$\lim_{d \rightarrow 0^+} \frac{\frac{d^k}{dz^k} \psi_d(z)}{d^j} = 0 \quad \forall j, k \in \mathbb{N}_0. \tag{4}$$

**Theorem 3** (Rough differential equation (RDE) solution). *In the compact regime, for each dataset  $X$  of size  $n$  (with entries  $x_i^\top$ ) and  $x \in S^{n_0-1}$ , there exists a sequence of paths  $v_{ij}^{(L)} : [0, 1] \rightarrow \mathbb{R}^n$  such that*

$$\lim_{L \rightarrow \infty} v_{ij}^{(L)}(t) = 0 \quad \forall t \in [0, 1]$$

*and the rough path lift  $\mathbf{v}^{(L)} : \Delta_{0,1} \rightarrow \mathbb{R}^{n \times n+1}$  of  $v_{ij}^{(L)}$  with  $p = 1$  **drives** the solution  $\mathbf{u}^{(L)}$  of a differential equation*

$$\frac{d}{dt} u_i^{(L)}(t) = 0 \quad \forall i \in \{1, \dots, n\},$$

*and whose projection  $(\mathbf{u}^{(L)})^1 = u^{(L)}$  onto 1-tensors satisfies the equality*

$$u_i^{(L)}(1) = \bar{\Theta}_\infty^{(L)}(x^\top X)^\top \left( \bar{\Theta}_\infty^{(L)}(X^\top X) \right)^{-1}.$$

Here,  $\Delta_{0,1}$  refers to the set  $\{(s, t) : 0 \leq s \leq t\}$ . Specifically,

$$\begin{aligned} \bar{\Theta}_\infty^{(L)}(x^\top X) (\bar{\Theta}_\infty(X^\top X))^{-1} &< C(x) \mathbf{1}_n^\top \\ \left\| (\bar{\Theta}(X^\top X))^{-1} \bar{\Theta}(X^\top x) \right\|_2 &\in \mathcal{O}(n) \end{aligned}$$

for  $L$  large enough. Moreover, when  $x$  is free and  $n$  is fixed, the function  $C$  is continuous and hence bounded on  $S^{n_0-1}$ .

*Proof.* We define the matrix  $A_n^{(L+1)}(t)$  with

$$\begin{aligned} A_n^{(L+1)}(t) &= \bar{\Theta}_\infty^{(L)}(XX^\top) + \psi_{\mathcal{D}}(2t-1) \left( \bar{\Theta}_\infty^{(L+1)}(XX^\top) - \bar{\Theta}_\infty^{(L)}(XX^\top) \right) \\ \mathcal{D} &= \det \left( \bar{\Theta}_\infty^{(L+1)}(XX^\top) \right) \det \left( \bar{\Theta}_\infty^{(L)}(XX^\top) \right) \end{aligned}$$

for  $t \in [0, 1]$  and dataset  $X = \{x_i\}_{i=1}^n$  of size  $n$ . We also define  $b_n^{(L+1)}(t)$  with

$$b_n^{(L+1)}(t) = \bar{\Theta}_\infty^{(L+1)}(x^\top X^\top).$$

From the system  $A_n^{(L+1)}(t)u(t) = b_n^{(L+1)}(t)$ , we take the derivative with respect to  $t$  and obtain the system

$$\left( \frac{d}{dt} A_n^{(L+1)}(t) \right) u(t) + A_n^{(L+1)}(t) \left( \frac{d}{dt} u(t) \right) = \frac{d}{dt} b_n^{(L+1)}(t).$$

Note that the solution  $u(t)$  depends implicitly on  $n, L$  and  $x$ . This will be made obvious later in the proof, but it is hidden for cleaner notation. By Cramer's rule, the solution to the system is

$$u'(t)_i = \frac{\sum_j \det \left( A_n^{(L+1)}(t) \leftarrow_{i,j} Z_A \right)}{\det \left( A_n^{(L+1)}(t) \right)} + \frac{\det \left( A_n^{(L+1)}(t) \leftarrow_{i,1} Z_b \right)}{\det \left( A_n^{(L+1)}(t) \right)}, \quad (5)$$

where  $Z_A, Z_b$  are defined by

$$Z_A = - \left( \frac{d}{dt} A_n^{(L+1)}(t) \right) \text{diag}(u(t)), \quad Z_b = \frac{d}{dt} b_n^{(L+1)}(t) = \mathbf{0}_n,$$

where the boldface  $\mathbf{0}_n$  denotes the vector of 0's of length  $n$ . By property (4) of  $\psi_{\mathcal{D}}$ , we obtain the sequence of inequalities

$$\begin{aligned} &\frac{\det \left( A_n^{(L+1)}(t) \leftarrow_{i,j} \frac{d}{dt} A_n^{(L+1)}(t) \right)}{\det \left( A_n^{(L+1)}(t) \right)} \quad (v_{(i,j)}) \\ &\leq \frac{\det \left( A_n^{(L+1)}(t) \leftarrow_{i,j} \frac{d}{dt} A_n^{(L+1)}(t) \right)}{\det \left( \bar{\Theta}_\infty^{(L+1)}(XX^\top) \right)^{\psi_{\mathcal{D}}(2t-1)} \det \left( \bar{\Theta}_\infty^{(L)}(XX^\top) \right)^{1-\psi_{\mathcal{D}}(2t-1)}} \\ &\leq \frac{\det \left( A_n^{(L+1)}(t) \leftarrow_{i,j} \frac{d}{dt} A_n^{(L+1)}(t) \right)}{\det \left( \bar{\Theta}_\infty^{(L+1)}(XX^\top) \right) \det \left( \bar{\Theta}_\infty^{(L)}(XX^\top) \right)} \rightarrow 0 \quad \text{as } L \rightarrow \infty, \end{aligned}$$

for  $L$  large enough. Note that we obtain the last inequality above through the fact that for  $L$  large enough, the strictly positive determinants are all smaller than 1. In addition, because, the function  $\psi_{\mathcal{D}}$  is infinitely smooth, the terms  $v_{(i,j)}$  are all of bounded total variation (see Definition 9 with  $p = 1$  in the appendix). For the same reason and using (4), we have that the  $v_{(i,j)}$  converge to 0 in the 1-variation metric. By Lyons Universal Limit theorem (Lyons, 1998) from rough path theory (see Definition 10 in the appendix), the solution  $u^{(L+1)}(t)$  (where we make the dependence on  $L + 1$  explicit) converges to the solution  $u_\infty(t)$  that solves the system  $u'_\infty(t) = \mathbf{0}_n$ . Hence,  $u_\infty(t)_i$  is a constant dependent on  $i$  and  $x$ . We

have a limiting solution  $u_\infty(t)$  that is bounded for each  $x$ . Because the Itô-Lyons map  $\Phi$  (Definition 12 in appendix) is continuous and locally Lipschitz in any  $p$ -norm variation topology, and because the rough path lift of  $v_{(i,j)}$  is continuous, the entire solution is bounded on the compact set  $S^{n_0-1}$ , i.e. there is a bound  $C' < \infty$  such that  $u_\infty(t) < C' \mathbf{1}_n$  for all  $x$  (note that  $u_\infty(t)$  depends on  $x$  but not  $C'$ ).  $\square$

An immediate consequence of this result is that if the depth  $L \in o(\min_{1 \leq l \leq L-1} n_l)$  and the width  $(\min_{1 \leq l \leq L-1} n_l)$  both go to infinity (i.e. the ratio of depth to width going to 0), the kernel of the output  $f_\infty$  of the neural network reaches a limiting expression. This limiting expression characterizes the effect of depth on infinitely wide fully-connected ReLU networks whose inputs lie on the sphere  $S^{n_0-1}$ . For ReLU networks, it is possible to easily extend this result to the non-compact regime (i.e. general domain  $\mathbb{R}^{n_0}$ ). By Proposition 2, there is a closed-form to  $\Theta_\infty$  for general data points in  $\mathbb{R}^{n_0}$ . In the statement of the proposition, the canonical projection on the sphere is provided, but a similar result is obtained for a stereographic projection.

## 6 EXPERIMENTS AND THEORETICAL IMPLICATIONS

In the proof of Theorem 3 from the previous section, we can identify the key properties used to derive the results, in order to distill the essence of the type of kernels that lead to similar limiting behaviour. We summarize these properties for data on  $S^{n_0-1}$ .

### Arbitrary sequence of kernels $\kappa^{(L)}$ satisfying the requirements of the theorem

- $\kappa^{(L)}(x, x) \geq \kappa^{(L)}(x_1, x_2)$  for any  $x, x_1, x_2 \in S^{n_0-1}$  and all  $L \in \mathbb{N}$ .
- There is some  $\hat{L} \in \mathbb{N}$  such that  $\kappa^{(L)}(XX^\top)$  is positive definite for any  $X = \{x_i \mid x_i \in S^{n_0-1}, i = 1, \dots, n\}$  of size  $n$  and  $L \geq \hat{L}$ .
- $\lim_{L \rightarrow \infty} \det(\bar{\kappa}^{(L)}(XX^\top)) = 0$  for data in  $S^{n_0-1}$ .

By inspection of the definition of  $\rho^{(L)}$ , it can be observed that it satisfies the criteria of the list above. Another example is given by the sequence  $\eta^{(L)}$  of kernels that are defined recursively by  $\eta^{(L+1)}(x, x') = h(\kappa^{(L)}(x, x'))$  and  $\eta^{(1)}(x, x') = x^\top x'$ , where  $h(z) = (1 + e^{-z})^{-2}$  (see Proposition 7 in the appendix).

In order to better understand the theoretical insights from the previous section, we empirically evaluate the convergence rates of  $\bar{\Theta}_\infty^{(L)}$ ,  $\rho^{(L)}$  and  $\eta$  as  $L$  increases. We illustrate this convergent behavior in figure 1, where we generate a dataset  $X$  and a point  $x$  from the uniform distribution ( $n_0 = 128$ ) and we canonically project them to the sphere. We then plot the evolution of the values for depths  $L = 1, \dots, 10$ . Each curve in the plots corresponds to a different pair of inputs (either from  $X$  or between  $x$  and  $X$ ; e.g.  $\bar{\Theta}_\infty^{(L)}(x_i, x_j)$ ). It is immediate at first glance that both  $\rho^{(L)}$  and  $\eta^{(L)}$  converge to a fixed value rather quickly as the depth increases. In contrast, from the plot of  $\bar{\Theta}_\infty^{(L)}(XX^\top)$ , while it is seemingly the case that off-diagonal values converge to some value strictly smaller than 1, Theorem 2 demonstrates that they converge to 1. The convergence rate is extremely slow and it is possible to observe this directly in the proof of the theorem:  $K$  is taken to be much larger than  $L$ , and given an approximation threshold of  $0 < \delta \ll 1$ , the exponential approximation to  $(1 - \delta)^{K+1}$  requires that  $K\delta \gg 0$ , i.e.  $K$  is much larger than  $\frac{1}{\delta}$ . The convergence rate to 1 is therefore extremely slow. While we do not explore any other kernels experimentally, we refer the reader to the list of criteria in this section to derive any other candidate to explore its convergence properties.

Note that Bietti & Bach (2020) and Li et al. (2024) tell us that the representation power of  $\Theta_\infty^{(L)}$  does not change as  $L \rightarrow \infty$ . However, if we apply our limiting result to the mean-field regime of Chizat & Bach (2018), we find that each particle approaches the deterministic limit by inspection of Proposition 3 and Theorem 3. It is therefore possible to analyze the many-particle limit of very wide and deep fully-connected neural networks since these are



well approximated by  $f_\tau$  for a proper stopping time  $\tau$  (Li et al., 2024). In addition, the proof technique of Theorem 3 can be adapted to other kernels that arise from other architectures such as CNNs.

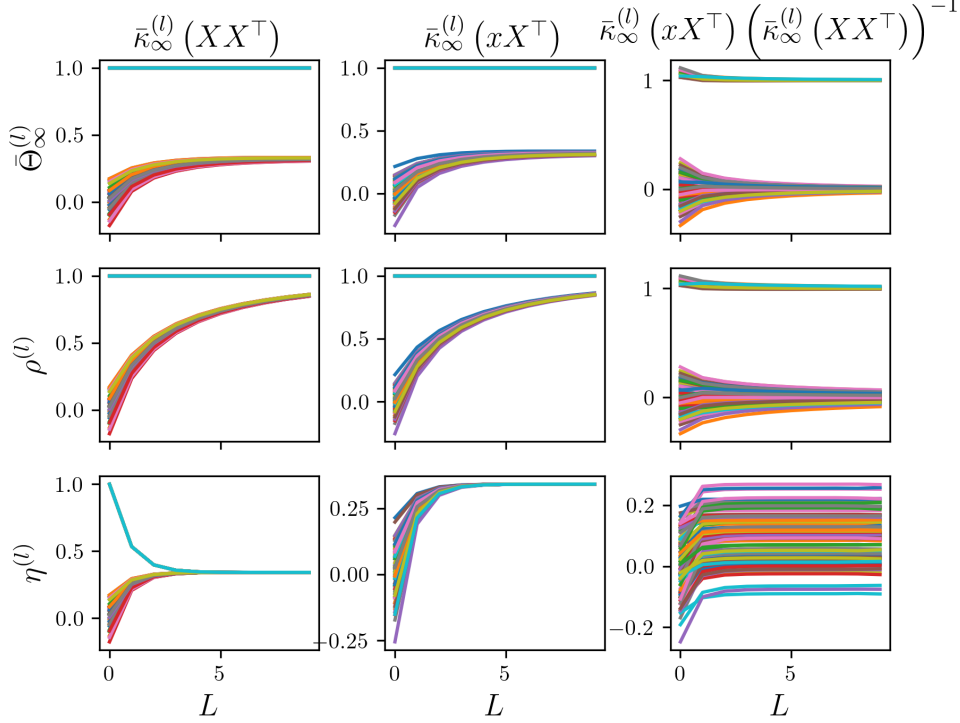


Figure 1: Convergence rate of  $\kappa$  on  $X$  and point  $x$ .

## 7 CONCLUSION

In this article, we provide a detailed analysis of behaviour of the deterministic kernel  $\Theta_\infty^{(L)}$  as  $L \rightarrow \infty$ . We observe that under the conditions of Jacot et al. (2018) and with  $L \in o(\min_{l=1,\dots,L-1} n_l)$ , a fully-connected ReLU neural network approaches a limiting solution for any  $x$  given a fixed dataset  $X$ . While our results are for data in  $S^{n_0-1}$ , we can extend this result to the general domain as long as there is no colinearity, or by taking the stereographic projection in a space with one additional dimension. We have shown the convergence rates of a non-exhaustive list of kernels as depth  $L$  increases. For the limiting kernels  $\kappa^{(L)} = \Theta_\infty^{(L)}$ , we have observed that the convergence is extremely slow and would require very large  $L$  before reaching the limit of  $\kappa_x \kappa^{-1}$ . Finally, we provided a list of key properties that were necessary to obtain our results to generalize to other kernels.

We believe that our work can help researchers to better understand the role of depth in the deterministic limiting kernels of overparameterized neural networks. Future research should envisage to better understand the behaviour of kernels that arise in the context of other architectures such as CNNs and architectures with skip connections. As mentioned, the setting of Hanin & Nica (2020) is outside the purview of our analysis as the stochasticity of the NTK becomes highly relevant when  $L \gg \max_{l=1,\dots,L-1} n_l$ . Nevertheless, by using the analytical tools of rough differential equations, we raise the hypothesis that there might exist a “pointwise” limit to the NTK when  $L \rightarrow \infty$ , where it is implied that the convergence is with respect to each sample path.

---

## 8 ACKNOWLEDGEMENTS

This work was funded by the NSERC Grant no. 2024-04051, NSERC Grant no. 2021-04378, and Canada CIFAR AI Chair. This research was enabled in part by support provided by Calcul Québec ([www.calculquebec.ca](http://www.calculquebec.ca)), Compute Canada ([www.computecanada.ca](http://www.computecanada.ca)) and Mila ([www.mila.quebec](http://www.mila.quebec)).

## REFERENCES

- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In *International conference on machine learning*, pp. 242–252. PMLR, 2019.
- Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International conference on machine learning*, pp. 322–332. PMLR, 2019a.
- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. *Advances in neural information processing systems*, 32, 2019b.
- Mikhail Belkin. Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation. *Acta Numerica*, 30:203–248, 2021.
- Alberto Bietti and Francis Bach. Deep equals shallow for relu networks in kernel regimes. *arXiv preprint arXiv:2009.14397*, 2020.
- Lenaic Chizat and Francis Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. *Advances in neural information processing systems*, 31, 2018.
- Youngmin Cho and Lawrence Saul. Kernel methods for deep learning. In Y. Bengio, D. Schuurmans, J. Lafferty, C. Williams, and A. Culotta (eds.), *Advances in Neural Information Processing Systems*, volume 22. Curran Associates, Inc., 2009. URL [https://proceedings.neurips.cc/paper\\_files/paper/2009/file/5751ec3e9a4feab575962e78e006250d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2009/file/5751ec3e9a4feab575962e78e006250d-Paper.pdf).
- Tilman Gneiting. Strictly and non-strictly positive definite functions on spheres. *Bernoulli*, 19(4):1327 – 1349, 2013. doi: 10.3150/12-BEJSP06. URL <https://doi.org/10.3150/12-BEJSP06>.
- Boris Hanin and Mihai Nica. Finite depth and width corrections to the neural tangent kernel. In *8th International Conference on Learning Representations, ICLR 2020*, 2020.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. *Advances in neural information processing systems*, 32, 2019.
- Jaehoon Lee, Samuel Schoenholz, Jeffrey Pennington, Ben Adlam, Lechao Xiao, Roman Novak, and Jascha Sohl-Dickstein. Finite versus infinite neural networks: an empirical study. *Advances in Neural Information Processing Systems*, 33:15156–15172, 2020.
- Yicheng Li, Zixiong Yu, Guhan Chen, and Qian Lin. On the eigenvalue decay rates of a class of neural-network related kernel functions defined on general domains. *Journal of Machine Learning Research*, 25(82):1–47, 2024.

- 
- Chaoyue Liu, Libin Zhu, and Misha Belkin. On the linearity of large non-linear models: when and why the tangent kernel is constant. *Advances in Neural Information Processing Systems*, 33:15954–15964, 2020.
- Chaoyue Liu, Libin Zhu, and Mikhail Belkin. Loss landscapes and optimization in over-parameterized non-linear systems and neural networks. *Applied and Computational Harmonic Analysis*, 59:85–116, 2022.
- Terry J Lyons. Differential equations driven by rough signals. *Revista Matemática Iberoamericana*, 14(2):215–310, 1998.
- Michael Murray, Hui Jin, Benjamin Bowman, and Guido Montufar. Characterizing the spectrum of the ntk via a power series expansion. In *International Conference on Learning Representations*, 2023.
- Quynh Nguyen, Marco Mondelli, and Guido F Montufar. Tight bounds on the smallest eigenvalue of the neural tangent kernel for deep relu networks. In *International Conference on Machine Learning*, pp. 8119–8129. PMLR, 2021.
- David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, Timothy Lillicrap, Karen Simonyan, and Demis Hassabis. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144, 2018. doi: 10.1126/science.aar6404. URL <https://www.science.org/doi/abs/10.1126/science.aar6404>.
- Russell Tsuchida, Fred Roosta, and Marcus Gallagher. Invariance of weight distributions in rectified mlps. In *International Conference on Machine Learning*, pp. 4995–5004. PMLR, 2018.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

## A THEORETICAL ASSUMPTIONS

For the proofs in Section 5, we make the following list of explicit assumptions:

**Activation  $\sigma$ :** ReLU activation.

**Output dimension:** The output dimension considered is  $n_L = 1$  for a network of depth  $L$ .

**Biases:** The neural networks do not contain any biases (i.e.  $\beta = 0$ ).

**Data:** The data has a fixed representation on  $S^{n_0-1} \subseteq \mathbb{R}^{n_0}$ .

**Kernel  $\Theta_\infty^{(L)}$ :** Limiting kernel of a fully-connected ReLU neural network of depth  $L$ , output dimension  $n_L = 1$  and without biases, under infinite width (i.e.  $\min_{1 \leq l \leq L-1} n_l \rightarrow \infty$ ).

Introducing biases via  $\beta > 0$  changes  $\Sigma^{(L)}$  and  $\Theta_\infty^{(L)}$ . Therefore, Proposition 1 does not hold in its exact form. However, we can show that in the limit  $L \rightarrow \infty$ , the kernel  $\bar{\Theta}_\infty^{(L)}$  converges to a constant matrix whose entries are all equal. By the same token, we can apply the ideas of Theorem 3, and we obtain a limiting expression for  $f_\infty$ . In the general list of properties for  $\kappa^{(L)}$  that is found in Section 6, the value of  $\hat{L} = 2$  as the limiting kernel becomes positive definite on the sphere for  $L \geq 2$ .

## B ADDITIONAL LEMMAS, PROPOSITIONS, AND THEOREMS

**Lemma 1** (Convergence of  $\rho^{(L)}$ ). *If  $\rho^{(1)}(x, x') \in ]-1, 1[$ , then  $\rho^{(L)}(x, x') \rightarrow 1$  as  $L \rightarrow \infty$ .*

*Proof.* We observe that  $\rho^{(L+1)}(x, x') = h(\rho^{(L)})$ , where  $h$  is the function

$$h(z) = \frac{z \arcsin(z)}{\pi} + \frac{\sqrt{1-z^2}}{\pi} + \frac{z}{2}.$$

This function has derivative  $h'(z) < 1$  on any fixed interval  $[a, b] \subseteq ]-1, 1[$  such that  $b \neq 1$  and is continuously differentiable on the same interval. Therefore, if  $H_n(z) = (h \circ \dots \circ h)(z)$  denotes the  $n^{\text{th}}$  power composition of  $h$ , we have  $H_L(z) \rightarrow \beta$ , some unique fixed-point, uniformly on  $] -1, 1[$ . This can be proved by observing that the domain of  $H_L$  is a compact set and that for any metric  $d$ , the distance  $d(H_L(z_1), H_L(z_2)) < d(z_1, z_2)$  for  $z_1 \neq z_2$  (first show that  $d(H_L(z), z)$  is continuous and has a minimum which is 0). Furthermore, we have that  $h(z) \geq z$  on  $[-1, 1]$ , with strict inequality on  $[-1, 1[$ . This implies that  $\beta = 1$  and the proof is finished.  $\square$

The following proposition is the early-stopping variant of proposition 3. We include a proof sketch for the reader, although we wish to underscore that this result is already known in the literature.

**Proposition 6** (closed-form for  $f_\tau$ ; section 5 from Jacot et al. (2018)). *Given  $\tau < \infty$  and the spectrum  $\Lambda(\kappa)$ , the output  $f_\tau$  is given by*

$$f_\tau(x) = f_0(x) + \kappa_x^\top \kappa^{-1} \bar{\Lambda}(y_0 - y^*),$$

where  $\kappa_x, y_0, y^*$  are as in proposition 3, and

$$\bar{\Lambda}_{ij} = \begin{cases} \exp(-\lambda_i \tau) - 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

is the diagonal matrix applying the function  $g_\tau(\lambda) = \exp(-\lambda \tau) - 1$  elementwise to the elements of  $\Lambda(\kappa)$ .

*Proof sketch.* From section 5 in Jacot et al. (2018), the output  $f_t(x)$  of the neural network is given by

$$f_t(x) = f_\infty(x) + \alpha_0^{(f_0 - f^*)}(x) + \sum_{i=1}^n \exp(-\lambda_i t) \alpha_i^{(f_0 - f^*)}(x)$$

where the terms  $\alpha_i^{(f_0-f^*)}$  are the component functions of the eigenvalue decomposition of  $f_0(x) - f^*(x)$ . Notice that we can write

$$\begin{aligned} f_\infty(x) + \alpha_0^{(f_0-f^*)}(x) &= f_0(x) - \sum_{i=1}^n \alpha_i^{(f_0-f^*)}(x) \\ &= f_0(x) - \kappa_x^\top \kappa^{-1}(y_0 - y^*), \end{aligned}$$

where the second equality follows from  $\alpha_i^{(f_0-f^*)}(x) = \kappa^{-1}(y_0 - y^*)_i e_i$  for standard basis vectors  $e_i \in \mathbb{R}^n$ . This implies the following equality:

$$f_t(x) = f_0(x) + \kappa_x^\top \kappa^{-1} \bar{\Lambda}(y_0 - y^*).$$

We have an expression for the output  $f_\tau(x)$  for any input  $x$  and stopping time  $\tau$ .  $\square$

The following proposition provides a proof sketch of the convergence of  $\eta^{(L)}$  as  $L \rightarrow \infty$  and other properties that satisfy the requirements for Theorem 3.

**Proposition 7** (Convergence of  $\eta^{(L)}$ ). *The values  $\eta^{(L)}(x, x')$  converge to a unique  $\beta > 0$  for all  $x, x' \in S^{n_0-1}$  as  $L \rightarrow \infty$ . Moreover, the kernels  $\eta^{(L)}$  are positive definite on  $S^{n_0-1}$  and satisfy  $\eta^{(L)}(x, x) \geq \eta^{(L)}(x_1, x_2)$ .*

*Proof sketch.* As  $L \rightarrow \infty$ , all values converge to the same limit since the derivative of  $h$  is strictly smaller than 1 on  $[-1, 1]$ . The kernels  $\eta^{(L)}$  also satisfy  $\eta^{(L)}(x, x) \geq \eta^{(L)}(x_1, x_2)$  since  $h(z)$  is monotone increasing in  $z \in [-1, 1]$ . The kernels are all positive definite on  $S^{n_0-1}$  since the function  $h$  is analytic on  $[-1, 1]$  and it has infinitely many even and odd terms in its power series expansion at 0 (i.e. not and even or odd function) that are strictly positive (Gneiting, 2013).  $\square$

The following is the proof of Proposition 4.

**Proposition 4** (Alternative formulation of  $\bar{\Theta}_\infty^{(L)}$ ). *The equality*

$$\bar{\Theta}_\infty^{(L+1)}(x, x') = \frac{L}{L+1} h'(\rho^{(L)}(x, x')) \bar{\Theta}_\infty^{(L)}(x, x') + \frac{1}{L+1} h(\rho^{(L)}(x, x'))$$

*holds  $\forall x, x' \in S^{n_0-1}$ . Moreover, the values in the normalized kernel are all found in the interval  $[0, 1]$ .*

*Proof.* The fact that the values are all contained in  $[0, 1]$  is immediate from the definition of  $\bar{\Theta}_\infty^{(L)}$  and Proposition 2. Now,

$$\begin{aligned} \Theta_\infty^{(L+1)}(x, x') &= \dot{\Sigma}^{(L+1)}(x, x') \Theta_\infty^{(L)}(x, x') + \Sigma^{(L+1)}(x, x') \\ &= \frac{1}{2} h'(\rho^{(L)}(x, x')) \Theta_\infty^{(L)}(x, x') \\ &\quad + \frac{1}{n_0 2^L} h(\rho^{(L)}(x, x')), \end{aligned}$$

where the first equality comes from Theorem 1 and the second equality uses Propositions 1 and 2. This implies

$$\begin{aligned} \bar{\Theta}_\infty^{(L+1)}(x, x') &= \frac{Ln_0 2^L}{(L+1)n_0 2^{L-1}} \\ &\quad \times \frac{n_0 2^{L-1} \Theta_\infty^{(L)}(x, x')}{L} \frac{1}{2} h'(\rho^{(L)}(x, x')) \\ &\quad + \frac{1}{L+1} h(\rho^{(L)}(x, x')) \\ &= \frac{L}{L+1} h'(\rho^{(L)}(x, x')) \bar{\Theta}_\infty^{(L)}(x, x') \\ &\quad + \frac{1}{L+1} h(\rho^{(L)}(x, x')), \end{aligned}$$

where the equalities come from the definition of  $\Theta_\infty^{(L+1)}(x, x')$  and the normalization factors of  $\bar{\Theta}_\infty^{(L+1)}$  and  $\bar{\Theta}_\infty^{(L)}$ .  $\square$

The following contains the proof of Theorem 2.

**Theorem 2** (Convergence of  $\bar{\Theta}_\infty^{(L)}$ ). *For any  $x, x' \in S^{n_0-1}$ , the value  $\bar{\Theta}_\infty^{(L)}(x, x')$  strictly increases to 1 as  $L \rightarrow \infty$ .*

*Proof.* We have a system describing  $\bar{\Theta}_\infty^{(L+1)}$ ,

$$\begin{pmatrix} \bar{\Theta}_\infty^{(L+1)} \\ 1 \end{pmatrix} = \begin{pmatrix} \frac{L}{L+1} h'(\rho^{(L)}(x, x')) & \frac{1}{L+1} h(\rho^{(L)}(x, x')) \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \bar{\Theta}_\infty^{(L)}(x, x') \\ 1 \end{pmatrix}$$

where we define the  $2 \times 2$  matrix on the right-hand side to be  $A^{(L)}$ . For now, fix  $L$  and observe the product

$$A^{(L+K)} \dots A^{(L)},$$

which we define to be the  $2 \times 2$  matrix  $A^{[L:L+K]}$ . First off, we can show that

$$A_{11}^{[L:L+K]} = \frac{L}{L+K+1} \prod_{k=0}^K h'(\rho^{(L+k)}(x, x')) \rightarrow 0$$

as  $K \rightarrow \infty$  (recall we fixed  $L$ ). The same convergence would be true if  $L \in o(K)$  as  $L, K \rightarrow \infty$ . We also obtain a sequence of inequalities bounding  $A_{12}^{[L:L+K]}$ , where

$$\begin{aligned} A_{12}^{[L:L+K]} &= \begin{pmatrix} A_{11}^{[L+2:L+K]} & A_{12}^{[L+2:L+K]} \end{pmatrix} \begin{pmatrix} \frac{1}{L+2} h'(\rho^{(L+1)}(x, x')) h(\rho^{(L)}(x, x')) + \frac{1}{L+2} h(\rho^{(L+1)}(x, x')) \\ 1 \end{pmatrix} \\ &= \frac{1}{L+K+1} \sum_{k=0}^K h(\rho^{(L+k)}(x, x')) \prod_{k'=k+1}^K h'(\rho^{(L+k')}(x, x')) \\ &\geq \frac{1-\delta}{L+K+1} \sum_{k=0}^K (1-\delta)^{K-k} \\ &= \frac{1-\delta}{L+K+1} \sum_{k=0}^K (1-\delta)^k \\ &= \frac{1-\delta}{L+K+1} \frac{1-(1-\delta)^{K+1}}{\delta} \end{aligned}$$

if  $L$  is large enough such that  $\min(h'(\rho^{(L)}(x, x')), h(\rho^{(L)}(x, x'))) \geq 1-\delta$  for a  $\delta \in ]0, 1[$ . The second equality was obtained by repeatedly doing the matrix-vector products using matrices  $A^{(L)}, \dots, A^{(L+K)}$ . We observe that by the identity  $(1+x)^\alpha \approx e^{\alpha x}$  for  $x$  small and  $\alpha x$  large, for  $\delta$  small enough

$$\begin{aligned} \frac{1-\delta}{L+K+1} \frac{1-(1-\delta)^{K+1}}{\delta} &\approx \frac{1-\delta}{L+K+1} \frac{1-e^{-K\delta}}{\delta} \\ &\approx (1-\delta) \frac{K}{L+K+1} \end{aligned}$$

which, once we unfreeze  $L$ , goes to  $1-\delta$  when  $L \in o(K)$  as  $L \rightarrow \infty$ . Since  $L$  is allowed to become arbitrarily large, we can take  $\delta$  arbitrarily small by Lemma 1 (note that  $K$  is still large enough such that  $K\delta \gg 0$ ).  $\square$

## C ROUGH PATH THEORY

In this section, we provide some background notions on rough path theory. For more information, see Lyons (1998).

**Definition 7** (Truncated tensor algebra  $T^{(m)}$ ). *The truncated tensor algebra  $T^{(m)}$  of  $\mathbb{R}^d$  is given by*

$$T^{(m)}(\mathbb{R}^d) = \bigoplus_{i=0}^m (\mathbb{R}^d)^{\otimes i},$$

where  $(\mathbb{R}^d)^{\otimes 0} \cong \mathbb{R}^d$  and  $\bigoplus$  is the direct sum operator.

**Definition 8** (Projective norm  $\pi$  on tensor powers). *Given any norm  $\|\cdot\|$  on  $\mathbb{R}^d$ , it can be used to define the projective norm  $\pi$  on any tensor product  $(\mathbb{R}^d)^{\otimes m}$  by*

$$\pi(x) = \inf \left\{ \sum_{i=1}^m \|a_i\| \|b_i\| : x = \sum_{i=1}^m a_i \otimes b_i \right\}.$$

**Definition 9** ( $p$ -variation metric). *Let  $\Delta_{0,1}$  be the simplex  $\{(s, t) : 0 \leq s \leq t \leq 1\}$  and  $p \geq 1$ . Let  $\mathbf{X}, \mathbf{Y}$  be continuous maps  $\Delta_{0,1} \rightarrow T^{(\lfloor p \rfloor)}(\mathbb{R}^d)$  and let  $\mathbf{X}^j$  (resp.  $\mathbf{Y}^j$ ) denote the projection of  $\mathbf{X}$  (resp.  $\mathbf{Y}$ ) onto its  $j$ -tensor component. The  $p$ -variation metric  $d_p$  is defined by*

$$d_p(\mathbf{X}, \mathbf{Y}) = \max_{j=1, \dots, \lfloor p \rfloor} \sup_{0=t_0 < t_1 < \dots < t_m=1} \left( \sum_{i=0}^{m-1} \|\mathbf{X}^j - \mathbf{Y}^j\|_{\frac{p}{j}}^{\frac{j}{p}} \right),$$

where the supremum is taken over all finite partitions  $\{0 = t_0 < t_1 < \dots < t_m = 1\}$  of  $[0, 1]$ .

**Definition 10** (Rough path (Lyons, 1998)). *A continuous function  $\mathbf{X} : \Delta_{0,1} \rightarrow T^{(\lfloor p \rfloor)}(\mathbb{R}^d)$  is a  $p$ -geometric rough path if there exists a sequence of paths with finite 1-variation (i.e. bounded variation)  $X(1), X(2), \dots$  such that*

$$\mathbf{X}(l)_{s,t} = \left( 1, \int_{s < s_1 < t} dX(l)_{s_1}, \dots, \int_{s < s_1 < \dots < s_{\lfloor p \rfloor} < t} dX(l)_{s_1} \otimes \dots \otimes dX(l)_{s_{\lfloor p \rfloor}} \right)$$

converges in the  $p$ -variation metric to  $\mathbf{X}$  as  $l \rightarrow \infty$ . The map from  $X(l)$  to  $\mathbf{X}(l)$  is called the **rough path lift** of  $X(l)$ .

**Definition 11** (Rough differential equation). *Let  $V_{i,j}$  for  $i = 1, \dots, e$  and  $j = 1, \dots, d$  be functions that have at least  $\lfloor p \rfloor$  bounded derivatives and the  $\lfloor p \rfloor$ -th derivatives are  $\alpha$ -Hölder continuous for  $\alpha > p - \lfloor p \rfloor$ . A rough differential equation takes the form*

$$(dY_t)_i = \sum_j^d V_{i,j}(Y_t) (dX_t)_j,$$

where  $Y$  is the solution to the differential equation,  $X$  is the driving signal, and both  $X$  and  $Y$  admit a rough path lift to a  $p$ -geometric rough path. If  $\mathbf{Y}$  and  $\mathbf{X}$  are the corresponding rough paths, we can also say that  $\mathbf{Y}$  solves the differential equation driven by  $\mathbf{X}$ .

**Definition 12** (Itô-Lyons map (Lyons, 1998)).  *$p$ -geometric rough paths  $\Delta_{0,1} \rightarrow T^{(\lfloor p \rfloor)}(\mathbb{R}^d)$  take value in the group  $G\Omega_p(\mathbb{R}^d)$  embedded in  $T^{(\lfloor p \rfloor)}(\mathbb{R}^d)$ . Given a rough differential equation, the Itô-Lyons map is the map  $\Phi : G\Omega_p(\mathbb{R}^d) \rightarrow G\Omega_p(\mathbb{R}^e)$  from a geometric rough path  $\mathbf{X}$  to a geometric rough path  $\mathbf{Y}$  solving the rough differential equation driven by  $\mathbf{X}$ .*

## D RELATED WORK SUMMARY

In this section, we present a summary of our results and their relation with the available literature of Section 2 on the subject. This summary is presented in Table 1.

---

Architecture	Activation	Results and relevant works
Wide, fixed depth, fully-connected	General	<ul style="list-style-type: none"> <li>• Convergence to NTK on the sphere and general domain; Hessian is approximately zero (Jacot et al., 2018; Belkin, 2021; Liu et al., 2020; 2022).</li> <li>• Spectrum characterized via Hermite expansion (Nguyen et al., 2021; Murray et al., 2023; Li et al., 2024).</li> <li>• Power law decay for NTK eigenvalues (Murray et al., 2023; Li et al., 2024).</li> <li>• Tight bound on smallest eigenvalue of NTK (Nguyen et al., 2021).</li> </ul>
Kernel model (uniform measure)	N/A	RKHS remains constant regardless of depth (Bietti & Bach, 2020).
Infinite width, depth, fully-connected	ReLU, general	<ul style="list-style-type: none"> <li>• NTK behaves stochastically (as <math>\lim_{L, n_1, \dots, n_{L-1} \rightarrow \infty} \max_{1 \leq l \leq L-1} \frac{L}{n_l} \gg 0</math>) (Hanin &amp; Nica, 2020).</li> <li>• Kernel limit as depth increases (as <math>\lim_{L, n_1, \dots, n_{L-1} \rightarrow \infty} \max_{1 \leq l \leq L-1} \frac{L}{n_l} = 0</math>) (<b>Our work 2025</b>).</li> <li>• Results for both spherical and general domains (Hanin &amp; Nica, 2020); (<b>Our work 2025</b>).</li> </ul>

---

Table 1: Summary of literature review and contributions