

# NOVEL KERNEL MODELS AND UNIFORM CONVERGENCE BOUNDS FOR NEURAL NETWORKS BEYOND THE OVER-PARAMETERIZED REGIME

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

This paper presents two models - called *global* and *local* models - of neural networks applicable to neural networks of arbitrary width, depth and topology, assuming only finite-energy neural activations. The first model is *exact* (unapproximated) and *global* (applicable for arbitrary weights), casting the neural network in reproducing kernel Banach space (RKBS). This leads to a width-independent (under usual scaling) bound on the Rademacher complexity of neural networks in terms of the spectral-norm of the weight matrices, which is depth-independent for sufficiently small weights. For illustrative purposes we consider how this bound may be applied to untrained networks with LeCun, He and Glorot initialization, discuss their connect to width and depth dependence in the Rademacher complexity bound, and suggest a modified He initialization that gives a depth-independent Rademacher complexity bound whp. The second model is exact and *local*, casting the *change* in neural network function resulting from a bounded change in weights and biases (ie. a training step) in reproducing kernel Hilbert space (RKHS) with a well-defined local-intrinsic neural kernel (LiNK). The neural tangent kernel (NTK) is shown to be a first-order approximation of the LiNK, so the local model gives insight into how the NTK model may be generalized outside of the over-parameterized limit. Analogous to the global model, a bound on the Rademacher complexity of network adaptation is obtained from the local model. Throughout the paper (a) dense feed-forward ReLU networks and (b) residual networks (ResNet) are used as illustrative examples and to provide insight into their operation and properties.

## 1 INTRODUCTION

The application of reproducing kernel Hilbert space (RKHS (Aronszajn, 1950)) and reproducing kernel Banach space (RKBS (Lin et al., 2022; Der & Lee, 2007; Zhang et al., 2009; Zhang & Zhang, 2012; Song et al., 2013; Sriperumbudur et al., 2011; Xu & Ye, 2014)) theory to the study of neural networks has a long history (Neal, 1996; Weinan et al., 2019; Parhi & Nowak, 2021; Lee et al., 2018; Matthews et al., 2018; Rahimi & Benjamin, 2009; Bach, 2014; 2017; Daniely et al., 2016; Daniely, 2017; Cho & Saul, 2009; Bartolucci et al., 2021; Spek et al., 2022). Neural tangent kernels (NTKs) are an exemplar of this approach, modeling training using a first-order (tangent) model. This approach has led to a wide body of work on convergence and generalization (Du et al., 2019b; Allen-Zhu et al., 2019; Du et al., 2019a; Zou et al., 2020; Zou & Gu, 2019; Arora et al., 2019b;a; Cao & Gu, 2019), mostly focused on the wide-network (over-parameterized) limit. In parallel, there is a growing body of work investigating the uniform convergence, complexity, and capacity of neural networks under various regimes (Neyshabur et al., 2015; 2018; 2019; 2017; Harvey et al., 2017; Bartlett et al., 2017; Golowich et al., 2018; Arora et al., 2018; Allen-Zhu et al., 2018; Dräxler et al., 2018; Li & Liang, 2018; Nagarajan & Kolter, 2019b; Zhou et al., 2019).

Nevertheless, as noted in eg. (Arora et al., 2019b; Lee et al., 2019; Bai & Lee, 2019), there are limits to the descriptive powers of NTK models. A gap has been observed between NTK-based predictions and actual performance (Arora et al., 2019b; Lee et al., 2019), and the validity of NTK models naturally breaks down outside of the over-parameterized limit, which has led to attempts to generalize NTK models outside of the over-parameterized - for example (Bell et al., 2023) use an

054 exact pathwise kernel, while (Shilton et al., 2023) presented an exact model for dense feedforward  
 055 neural networks with smooth activations in RKBS, (Bartolucci et al., 2023; Sanders, 2020; Parhi &  
 056 Nowak, 2021; Unser, 2021; 2019) have explored links to RKBS, and (Bai & Lee, 2019) explored  
 057 higher-order approximations. However the assumptions made (smoothness, over-parameterization  
 058 etc) and approximations used limit application and raise the question of whether it is possible to  
 059 instead formulate *exact, non-approximate* models for neural networks that may be used for similar  
 060 ends. With regard to uniform convergence, some authors such as (Nagarajan & Kolter, 2019a) have  
 061 argued that such methods may be unable to explain generalization for neural networks at all due to the  
 062 behavior of the Rademacher complexity (Bartlett & Mendelson, 2002; Steinwart & Christman, 2008).

063 In this paper we simultaneously address two questions, namely: (1) is it possible to formulate an  
 064 exact (non-approximate) model for a wide class of neural networks, thereby avoiding entirely the  
 065 question of gaps between real performance and model prediction; and (2) can such a model be used  
 066 to derive general, non-vacuous, training-independent bounds on uniform convergence. We address  
 067 both questions with the following contributions:

- 068 1. **Exact global model:** for arbitrary neural network topologies with arbitrary weights and  
 069 biases and finite-energy activations, we construct a model that recasts neural networks and  
 070 bilinear products in a reproducing kernel Banach space (RKBS). This leads to:
  - 071 (a) **Rademacher Complexity Bound:** using the global model, we bound the Rademacher  
 072 complexity as a function of the spectral-norms of the weight matrices. We show that  
 073 this bound is *width-independent* for standard weight-scaling, and depth-independent in  
 074 the unbiased case for small weight matrices if the neural activations are  $L$ -Lipschitz.  
 075 For networks satisfying this constraints, we prove that the Rademacher complexity is  
 076 bounded as  $\mathcal{R}_N(\mathcal{F}) \leq \frac{1}{\sqrt{N}}$  for training set size  $N$ .
  - 077 (b) **Rademacher Complexity of Randomly Initialized Networks:** we discuss the  
 078 Rademacher complexity of randomly initialized networks, in particular following Le-  
 079 Cun, He and Glorot initialization. We analyse the width- and depth- dependence of the  
 080 Rademacher complexity bound for these initializations and present modified He and  
 081 Glorot initializations for which  $\mathcal{R}_N(\mathcal{F}) \leq \frac{1}{\sqrt{N}}$  with high probability.
- 082 2. **Exact local model:** again for arbitrary neural network topologies with arbitrary weights  
 083 and biases and finite-energy activations, we construct a model that recasts the *change* in  
 084 neural network operation due to a (spectral-norm) bounded change in weights and biases in  
 085 a reproducing kernel Hilbert space (RKHS), with a locally-intrinsic neural kernel (LiNK)  
 086 defined by the network topology, neural activations and initial weights. This leads to:
  - 087 (a) **Rademacher Complexity Bound:** analogous to the global model, the local model  
 088 leads to a bound on the Rademacher complexity as a function of the spectral-norms of  
 089 the *change* to the weight matrices.
  - 090 (b) **Connection with NTK:** we prove that the NTK is a first-order approximation LiNK  
 091 in our local model, casting light on higher-order generalization of NTK models.

092 The paper is organized as follows. We first discuss the underlying scope and setting of the paper,  
 093 and the relevant notions and notations used (section 2), as well as related work (section 3). The  
 094 necessary mathematical background on Hermite polynomials is provided in section 4, including some  
 095 discussion regarding how this will inform our contribution. We present our global model in section  
 096 5 before proceeding to use this model to derive bounds on Rademacher complexity in section 5.2.  
 097 We finish by presenting our local model in section 6, which is applied to the problem of uniform  
 098 convergence in section 6.2.

## 099 1.1 MATHEMATICAL NOTATION AND INDEXING CONVENTIONS

101 We use  $\mathbb{N} = \{0, 1, 2, \dots\}$ ,  $\mathbb{N}_n = \{0, 1, 2, \dots, n - 1\}$ ,  $\mathbb{Z}_+ = \{1, 2, 3, \dots\}$ .  $|\mathbb{A}|$  is the number of  
 102 elements in set  $\mathbb{A}$ .  $\text{Span}(\mathcal{X})$  is the linear span of  $\mathcal{X}$ .  $[a]_+ = \max\{a, 0\}$ .  $f^{(n)}$  is the  $n^{\text{th}}$  derivative of  
 103  $f$ .  $L^2(\mathbb{R}, e^{-x^2}) = \{\tau : \mathbb{R} \rightarrow \mathbb{R} \mid \int_{-\infty}^{\infty} |\tau(\zeta)|^2 e^{-\zeta^2} d\zeta < \infty\}$  is the set of finite-energy functions.  $He_k$   
 104 are the (probabilist’s) Hermite polynomials ( $k \in \mathbb{N}$ ).  $He_k = He_k(0)$  are the Hermite numbers.

105 Column vectors are denoted  $\mathbf{a}, \mathbf{b}, \dots$ , with elements  $a_i, b_j, \dots$   $|\mathbf{a}|$  and  $\text{sgn}(\mathbf{a})$  are the element-wise  
 106 absolute and sign.  $\|\mathbf{a}\|_2 = (\sum_i |a_i|^2)^{1/2}$  is the Euclidean norm. Matrices are denoted  $\mathbf{A}, \mathbf{B}, \dots$ ,  
 107 with elements  $A_{i,i'}$ , rows  $\mathbf{A}_i$ , and columns  $\mathbf{A}_{i'}$ .  $\mathbf{A} \odot \mathbf{B}$  is the Hadamard product.  $\mathbf{A} \otimes \mathbf{B}$  is the

108 Kronecker product.  $\mathbf{A}^{\otimes k} = \mathbf{A} \otimes \mathbf{A} \otimes \dots \otimes \mathbf{A}$  is the Kronecker power.  $\text{diag}_i(\mathbf{A}_i)$  is a block-  
 109 diagonal matrix with diagonal blocks  $\mathbf{A}_i$ .  $\|\mathbf{A}\|_2 = \sup_{\mathbf{x}: \|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2$  is the spectral norm.

110  $\langle \cdot, \cdot \rangle$  denotes an inner-product, where  $\langle \mathbf{a}, \mathbf{a}' \rangle_{\mathbf{g}} = \sum_i g_i a_i a'_i$  for metric  $\mathbf{g} > \mathbf{0}$ , and unless otherwise  
 111 stated  $\langle \mathbf{a}, \mathbf{a}' \rangle = \langle \mathbf{a}, \mathbf{a}' \rangle_{\mathbf{1}}$ .  $\langle \cdot, \cdot \rangle$  denotes a bi-linear product, where  $\langle \mathbf{a}, \mathbf{a}' \rangle_{\mathbf{g}} = \sum_i g_i a_i a'_i$  for (indefinite)  
 112 metric  $\mathbf{g}$ , and we follow the convention  $\langle \mathbf{A}, \mathbf{a}' \rangle_{\mathbf{g}} = [\langle \mathbf{A}_{:,j}, \mathbf{a}' \rangle_{\mathbf{g}}]_j$ ,  $\langle \mathbf{A}, \mathbf{a}' \rangle_{\mathbf{G}} = [\langle \mathbf{A}_{:,j}, \mathbf{a}' \rangle_{\mathbf{G}_{:,j}}]_j$ .

113 **Indexing:**  $\tilde{j}, j \in \mathbb{N}_D$  index nodes in the computational graph.  $\tilde{v}_j \in \mathbb{N}_{\tilde{H}^{[j]}}$  and  $i_j \in \mathbb{N}_{H^{[j]}}$  index  
 114 inputs and outputs to node  $j$ , respectively, where  $\tilde{H}^{[j]}$  and  $H^{[j]}$  are the fan-in and fan-out of that  
 115 node.  $\mathbb{P}^{[j]}$  is the antecedent set of node  $j$  (set of nodes feeding into node  $j$ ), and  $p^{[j]} = |\mathbb{P}^{[j]}|$  the  
 116 node’s in-degree. Node  $D - 1$  is the output node, and  $j = -1$  indicates the (virtual) input node.  
 117  
 118  
 119

## 120 2 SETTING AND ASSUMPTIONS

121  
 122 In this paper we will deal with neural networks defined by directed acyclic graphs (aka computational  
 123 skeletons) as per (Daniely et al., 2016), as this description is flexible enough to allow us to study both  
 124 simple, fully-connected feed-forward networks such as ReLU and also non-trivial topologies such  
 125 residual networks (ResNet). In this scheme networks are characterized by nodes and edges, ie.:

- 126 1. **Nodes:** a set of  $D$  nodes, indexed by  $j \in \mathbb{N}_D$ , where node  $j = D - 1$  is the output node  
 127 and we include a virtual input node indexed as  $j = -1$ . A described shortly, a node  $j$  is  
 128 characterized by a weight matrix  $\mathbf{W}^{[j]} \in \mathbb{R}^{\tilde{H}^{[j]} \times H^{[j]}}$  (sometimes split into individual sub-  
 129 matrices  $\mathbf{W}^{[\tilde{j},j]} \in \mathbb{R}^{H^{[\tilde{j}]} \times H^{[j]}}$  per incoming edge) and bias vector  $\mathbf{b}^{[j]} \in \mathbb{R}^{H^{[j]}}$ .
- 130 2. **Edges:** a set of directed edges  $(\tilde{j} \rightarrow j) \in (\mathbb{N}_D \cup \{-1\}) \times \mathbb{N}_D$  characterized by neural  
 131 activation functions  $\tau^{[\tilde{j},j]} : \mathbb{R} \rightarrow \mathbb{R}$ , joining nodes to form a directed acyclic graph (DAG).  
 132

133 Given an input  $\mathbf{x}$ , data flows from the virtual input node  $j = -1$ , along the the edges and through the  
 134 nodes of the DAG to the output node  $j = D - 1$ , as defined by the recursive equation:

$$135 \mathbf{f}(\mathbf{x}; \Theta) = \mathbf{x}^{[D-1]} \left\{ \begin{array}{l} \tilde{\mathbf{x}}^{[j]} = [\tilde{\mathbf{x}}^{[\tilde{j},j]}]_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}}, \tilde{\mathbf{x}}^{[\tilde{j},j]} = \tau^{[\tilde{j},j]}(\mathbf{x}^{[\tilde{j}]}) \\ \mathbf{x}^{[j]} = \mathbf{W}^{[j]\text{T}} \tilde{\mathbf{x}}^{[j]} + \gamma \mathbf{b}^{[j]} \\ \left( = \sum_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \mathbf{W}^{[\tilde{j},j]\text{T}} \tilde{\mathbf{x}}^{[\tilde{j},j]} + \gamma \mathbf{b}^{[j]} \right) \end{array} \right\} \begin{array}{l} \forall j \in \mathbb{N}_D, \tilde{j} \in \tilde{\mathbb{P}}^{[j]} \\ \mathbf{x}^{[-1]} = \mathbf{x} \end{array} \quad (1)$$

136 For node  $j$  we define the antecedent set  $\tilde{\mathbb{P}}^{[j]} \subset \mathbb{N}_D \cup \{-1\}$ , so the node has in-degree  $p^{[j]} = |\tilde{\mathbb{P}}^{[j]}|$ ,  
 137 fan-out  $H^{[j]}$  (with input dimension  $H^{[-1]} = n$  and output dimension  $H^{[D-1]} = m$ ) and fan-in  
 138  $\tilde{H}^{[j]} = \sum_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} H^{[\tilde{j}]}$ . In constructing our models we assume:  
 139

- 140 1. **Bounded inputs:**  $\|\mathbf{x}\|_2 \leq 1$ .
- 141 2. **Finite weights and biases:**  $\|\mathbf{W}^{[j]}\|_2, \|\mathbf{b}^{[j]}\|_2 < \infty$ .
- 142 3. **Finite activations:**  $\tau^{[\tilde{j},j]} \in L^2(\mathbb{R}, e^{-\zeta^2}) = \{\tau : \mathbb{R} \rightarrow \mathbb{R} \mid \int_{-\infty}^{\infty} |\tau(\zeta)|^2 e^{-\zeta^2} d\zeta < \infty\}$ .

143 We let  $\Theta = \{\mathbf{W}^{[j]}, \mathbf{b}^{[j]} : j \in \mathbb{N}_D\}$  denote the collection of all weights and biases. We also denote the  
 144 set of inputs satisfying the bounded input assumption  $\mathbb{X}$ , and the set of weights and biases satisfying  
 145 the finiteness assumption  $\mathbb{W}$ , so  $\mathbf{f} : \mathbb{X} \times \mathbb{W} \rightarrow \mathbb{R}^m$ . Given a training set  $\mathcal{D} = \{(\mathbf{x}_{\{i\}}, \mathbf{y}_{\{i\}}) \in \mathbb{R}^n \times$   
 146  $\mathbb{R}^m : i \in \mathbb{N}_N\}$  we assume the goal is to minimize the risk (for loss  $L$ , regularizer  $r$ , trade-off  $\lambda \in \mathbb{R}_+$ ):  
 147

$$148 \Theta^* = \underset{\Theta \in \mathbb{W}}{\text{argmin}} \sum_i L(\mathbf{f}(\mathbf{x}_{\{i\}}; \Theta) - \mathbf{y}_{\{i\}}) + \lambda r(\Theta) \quad (2)$$

149 When constructing and analysing the network for the global model, for all  $j \in \mathbb{N}_D, \tilde{j} \in \tilde{\mathbb{P}}^{[j]}$  we find  
 150 it convenient to define a nominal upper bounds  $\mu^{[\tilde{j},j]}$  on the spectral norm of the weight matrices and  
 151  $\beta^{[j]}$  on the Euclidean norm of the biases:  
 152

$$153 \|\mathbf{W}^{[\tilde{j},j]}\|_2 \leq \mu^{[\tilde{j},j]} \text{ and } \|\mathbf{b}^{[j]}\|_2 \leq \beta^{[j]} \quad (3)$$

154 Similarly for the local model (and in the analysis of the global model) we find it convenient to define  
 155 a nominal bound  $\mu^{[j]}$  so that  $\|\mathbf{W}^{[j]}\|_2 \leq \mu^{[j]}$ . Building these into the model helps us to simplify  
 156 our results later. *It is important to note that these are convenience factors only and may be as large*  
 157 *as necessary to satisfy (3). The only restriction we make here is that  $\mu^{[\tilde{j},j]}$  and  $\beta^{[j]}$  must be finite.*  
 158  
 159  
 160  
 161

For example for randomly-initialized, untrained neural networks we may derive appropriate, high-probability upper bounds  $\mu^{[j]}$  and  $\beta^{[j]}$  by considering the distribution from which the weights and biases are drawn. In the usual case  $W_{\tilde{i}_j, i_j}^{[j]} \sim \mathcal{N}(0, \sigma^{[j]2})$  we have that  $\|\mathbf{W}_{:\tilde{i}_j}^{[j,j]}\|_2^2 \sim \sigma^{[j]2} \chi_{H^{[j]}}^2$ , so:

$$\|\mathbf{W}_{:\tilde{i}_j}^{[j,j]}\|_2^2 \leq \sigma^{[j]2} \left( H^{[j]} + 2\sqrt{H^{[j]}} \ln\left(\frac{DH^{[j]}}{2\epsilon}\right) + 2 \ln\left(\frac{DH^{[j]}}{\epsilon}\right) \right) \quad (4)$$

whp  $\geq 1 - \epsilon$  simultaneously  $\forall j, i_j$  (see eg. (Laurent & Massart, 2000, Lemma 1, pg 1325)). From this we may derive bounds for several standard initialization schemes, for example:

1. **LeCun** ( $\sigma^{[j]2} = \frac{1}{H^{[j]}}$ ):  $\mu^{[\tilde{j},j]2} = \frac{H^{[j]}}{H^{[j]}} + \frac{2\sqrt{H^{[j]}}}{H^{[j]}} \ln\left(\frac{DH^{[j]}}{2\epsilon}\right) + \frac{2}{H^{[j]}} \ln\left(\frac{DH^{[j]}}{\epsilon}\right)$ .
2. **He** ( $\sigma^{[j]2} = \frac{1}{H^{[j]}}$ ):  $\mu^{[\tilde{j},j]2} = \frac{H^{[j]}}{H^{[j]}} + \frac{2\sqrt{H^{[j]}}}{H^{[j]}} \ln\left(\frac{DH^{[j]}}{2\epsilon}\right) + \frac{2}{H^{[j]}} \ln\left(\frac{DH^{[j]}}{\epsilon}\right)$ .
3. **Glorot** ( $\sigma^{[j]2} = \frac{1}{H^{[j]} + \tilde{H}^{[j]}}$ ):  $\mu^{[\tilde{j},j]2} = \frac{H^{[j]}}{H^{[j]} + \tilde{H}^{[j]}} + \frac{2\sqrt{H^{[j]}}}{H^{[j]} + \tilde{H}^{[j]}} \ln\left(\frac{DH^{[j]}}{2\epsilon}\right) + \frac{2}{H^{[j]} + \tilde{H}^{[j]}} \ln\left(\frac{DH^{[j]}}{\epsilon}\right)$ .

To illustrate our results we use the following network topologies (Glorot et al., 2011; He et al., 2016):

1. **Feedforward ReLU**: fully connected, unbiased, layerwise, feedforward, ReLU activations:  $\tilde{\mathbb{P}}^{[j]} = \{j-1\}$ ,  $\gamma = 0 \forall j$ ;  $\tau^{[j-1,j]}(\zeta) = [\zeta]_+ \forall j \neq 0, D-1$ ;  $\tau^{[-1,0]}(\zeta) = \tau^{[D-2,D-1]}(\zeta) = \zeta$ .
2. **Residual Network (ResNet)**: unbiased, alternating ReLU/skip network with  $D \in 2\mathbb{Z}_+$ :  $\tilde{\mathbb{P}}^{[j]} = \{j-1\} \forall j$  even;  $\tilde{\mathbb{P}}^{[j]} = \{j-1, j-2\}$ ,  $\mathbf{W}^{[j]} = [\mathbf{W}_C^{[j]}; \frac{1}{2}\mathbf{I}] \forall j$  odd;  $\gamma = 0 \forall j$ ;  $\tau^{[j-1,j]}(\zeta) = [\zeta]_+ \forall j \neq 0, D-1$ ;  $\tau^{[j-2,j]}(\zeta) = \zeta \forall j$  odd;  $\tau^{[-1,0]}(\zeta) = \tau^{[D-2,D-1]}(\zeta) = \zeta$ .

### 3 RELATED WORK

The use of kernel methods to model neural networks dates to at least (Neal, 1996). Fixing all weights and biases except for node  $j$  and defining feature map  $\varphi^{[j]}(\mathbf{x}) = [\tilde{\mathbf{x}}^{[j]}; \gamma]$ , the network can be written as  $\mathbf{f}(\mathbf{x}; \Theta) = \mathbf{q}^{[j]}([\mathbf{W}^{[j]\top}; \mathbf{b}^{[j]})\varphi^{[j]}(\mathbf{x}, \mathbf{x})$  for fixed  $\mathbf{q}^{[j]}$ , and (2) becomes kernel regression with NNGP (neural-network Gaussian process) kernel:

$$\begin{aligned} \Sigma^{[j]}(\mathbf{x}, \mathbf{x}') &= \mathbb{E}_k \left[ \varphi_k^{[j]}(\mathbf{x}) \varphi_k^{[j]}(\mathbf{x}') \right] = \mathbb{E}_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \left[ \Sigma^{[\tilde{j},j]}(\mathbf{x}, \mathbf{x}') \right] \\ \Sigma^{[\tilde{j},j]}(\mathbf{x}, \mathbf{x}') &= \gamma^2 + \mathbb{E}_{i_{\tilde{j}}} \left[ \tilde{x}_{i_{\tilde{j}}}^{[\tilde{j},j]} \tilde{x}'_{i_{\tilde{j}}} \right] \end{aligned} \quad (5)$$

In the wide limit, for suitable initialization, (5) is deterministic (dependent on the distribution  $\Theta \sim \nu$ ), and it can be demonstrated that  $\mathbf{x}^{[j]}(\cdot) \sim \text{GP}(0, \Sigma^{[j]})$  (Neal, 1996; Lee et al., 2018; Matthews et al., 2018; Garriga-Alonso et al., 2019; Novak et al., 2019). This is the NNGP model, and may be used to eg. derive insights into the types of function the network is best suited to model. The NNGP kernel for our ReLU example is the arc-cosine kernel (Cho & Saul, 2009).

To model training, neural tangent kernels (NTKs (Jacot et al., 2018; Arora et al., 2019b)) form the basis of a first-order approximation of the behavior of a neural network as the weights and biases vary about their initialization  $\Theta$ , i.e.  $\mathbf{f}(\mathbf{x}; \Theta + \Delta\Theta) \approx \mathbf{f}(\mathbf{x}; \Theta) + \Delta\Theta^\top \nabla_{\Theta} \mathbf{f}(\mathbf{x}; \Theta)$ . Training is cast as kernel regression using the neural tangent kernel (NTK), recursively defined as:

$$\begin{aligned} K_{\text{NTK}}(\mathbf{x}, \mathbf{x}') &= \mathbb{E}_k \left[ \nabla_{\Theta_k} \mathbf{f}(\mathbf{x}; \Theta)^\top \nabla_{\Theta_k} \mathbf{f}(\mathbf{x}'; \Theta) \right] = K_{\text{NTK}}^{[D-1]}(\mathbf{x}, \mathbf{x}') \\ K_{\text{NTK}}^{[j]}(\mathbf{x}, \mathbf{x}') &= \Sigma^{[j]}(\mathbf{x}, \mathbf{x}') + \mathbb{E}_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \left[ \theta^{[\tilde{j},j]}(\mathbf{x}, \mathbf{x}') K_{\text{NTK}}^{[\tilde{j}]}(\mathbf{x}, \mathbf{x}') \right] \\ \theta^{[\tilde{j},j]}(\mathbf{x}, \mathbf{x}') &= \mathbb{E}_{i_{\tilde{j}}} \left[ \left\langle \tau^{[\tilde{j},j](1)}(x_{i_{\tilde{j}}}^{[\tilde{j}]}) \tau^{[\tilde{j},j](1)}(x'_{i_{\tilde{j}}} \right) \right\rangle \forall \tilde{j} \in \tilde{\mathbb{P}}^{[j]} \right] \end{aligned} \quad (6)$$

and  $K_{\text{NTK}}^{[-1]}(\mathbf{x}, \mathbf{x}') = 0$ . The NTK model is accurate for small variations in weights/biases (the lazy regime). In the infinitely wide limit the NTK is deterministic, and weights/biases remain close to their initial values, leading to the gradient flow model where weights flow rather than change in discrete steps during training. This approach gives insight in areas including convergence and generalization (Du et al., 2019b; Allen-Zhu et al., 2019; Du et al., 2019a; Zou et al., 2020; Zou & Gu, 2019; Arora et al., 2019b;a; Cao & Gu, 2019).

Beyond this first-order model, while NTKs have made significant progress, a gap has been observed

	Incoming Edge Feature Map	Node Feature Map
216		
217		
218		
219		
220	Feature Maps	
221	$\tilde{\phi}^{[\bar{j},j]}(\mathbf{x}) = \frac{1}{\phi} \begin{bmatrix} a_{(0)k}^{[\bar{j},j]} \left[ \binom{k}{l} \frac{1}{2} \left( \sqrt{\frac{1}{\bar{s}^{[\bar{j},j]-1}} (\phi^2)} \phi^{[l]}(\mathbf{x}) \right)^{\otimes l} \right]_{1 \leq l \leq k, k \geq 1} \\ a_{(0)k}^{[\bar{j},j]} \left[ \binom{k}{l} \frac{1}{2} \left( \frac{1}{\sqrt{\frac{1}{\bar{s}^{[\bar{j},j]-1}} (\phi^2)}} \Psi^{[l]}(\Theta) \right)^{\otimes l} \right]_{1 \leq l \leq k, k \geq 1} \end{bmatrix}$	$\phi^{[j]}(\mathbf{x}) = \frac{1}{\sqrt{\gamma^2+1}} \begin{bmatrix} \left[ \sqrt{\frac{H^{[j]}}{H^{[j]}}} \tilde{\phi}^{[\bar{j},j]}(\mathbf{x}) \right]_{j \in \bar{p}^{[j]}} \\ \mathbf{b}^{[j]T} + \sum_{j \in \bar{p}^{[j]}} \frac{\tau^{[\bar{j},j]}(0)}{\gamma} \mathbf{1}_{H^{[j]}}^T \mathbf{W}^{[\bar{j},j]} \end{bmatrix}$
222	$\tilde{\Psi}^{[\bar{j},j]}(\Theta) = \tilde{\phi} \begin{bmatrix} a_{(0)k}^{[\bar{j},j]} \left[ \binom{k}{l} \frac{1}{2} \left( \frac{1}{\sqrt{\frac{1}{\bar{s}^{[\bar{j},j]-1}} (\phi^2)}} \Psi^{[l]}(\Theta) \right)^{\otimes l} \right]_{1 \leq l \leq k, k \geq 1} \\ a_{(0)k}^{[\bar{j},j]} \left[ \binom{k}{l} \frac{1}{2} \left( \frac{1}{\sqrt{\frac{1}{\bar{s}^{[\bar{j},j]-1}} (\phi^2)}} \Psi^{[l]}(\Theta) \right)^{\otimes l} \right]_{1 \leq l \leq k, k \geq 1} \end{bmatrix}$	$\Psi^{[j]}(\Theta) = \sqrt{\gamma^2+1} \begin{bmatrix} \text{diag}_{j \in \bar{p}^{[j]}} \left( \sqrt{\frac{H^{[j]}}{H^{[j]}}} \tilde{\phi} \tilde{\Psi}^{[\bar{j},j]}(\Theta) \mathbf{W}^{[\bar{j},j]} \right) \\ \mathbf{g}^{[j]} \end{bmatrix}$
223		
224	Norm Bounds	
225	$\ \tilde{\phi}^{[\bar{j},j]}(\mathbf{x})\ _2^2 \in \left[ \phi_+^{[\bar{j},j]2} = \frac{1}{\phi^2} \bar{s}^{[\bar{j},j]} \left( \bar{s}^{[\bar{j},j]-1} (\phi^2) \phi_+^{[l]2} \right), \phi_-^{[\bar{j},j]2} = 1 \right]$	$\ \phi^{[j]}(\mathbf{x})\ _2^2 \in \left[ \phi_+^{[j]2} = \frac{1}{\gamma^2+1} \left( \gamma^2 + \sum_{j \in \bar{p}^{[j]}} \frac{H^{[j]}}{H^{[j]}} \tilde{\phi}_+^{[\bar{j},j]2} \right), \phi_-^{[j]2} = 1 \right]$
226	$\ \tilde{\Psi}^{[\bar{j},j]}(\Theta)\ _2^2 \leq \tilde{\psi}^{[\bar{j},j]2} = \tilde{\phi}^2 \bar{s}^{[\bar{j},j]} \left( \frac{1}{\bar{s}^{[\bar{j},j]-1} (\phi^2)} \psi^{[l]2} \right)$	$\ \Psi^{[j]}(\Theta)\ _2^2 \leq \psi^{[j]2} = (\gamma^2+1) \left( \left( \beta^{[j]} + \frac{\mu^{[j]}  \tau^{[\bar{j},j]}(0) }{\gamma} \right)^2 + \sum_{j \in \bar{p}^{[j]}} \frac{H^{[j]}}{H^{[j]}} \tilde{\psi}^{[\bar{j},j]2} \mu^{[\bar{j},j]2} \right)$
227	$\ \tilde{\Psi}^{[\bar{j},j]}(\Theta)\ _{\text{He}^{[\bar{j}]}}^2 \leq \tilde{\psi}^{[\bar{j},j]2} = \sup_{\substack{\phi_+^{[l]} \leq \phi_+^{[l]} \leq 1 \\ -\phi_-^{[l]} \leq \phi_-^{[l]} \leq \phi_-^{[l]}}} \left\{ \frac{\phi_+^{2, [\bar{j},j]} (\phi_+^{[l]} \psi^{[l]})^2}{\bar{s}^{[\bar{j},j]} (\bar{s}^{[\bar{j},j]-1} (\phi^2) \phi_+^{[l]2})} \right\}$	$\ \Psi^{[j]}(\Theta)\ _{\text{He}^{[j]}}^2 \leq \psi^{[j]2} = (\gamma^2+1) \left( \left( \beta^{[j]} + \frac{\mu^{[j]}  \tau^{[\bar{j},j]}(0) }{\gamma} \right)^2 + \sum_{j \in \bar{p}^{[j]}} \frac{H^{[j]}}{H^{[j]}} \psi^{[\bar{j},j]2} \mu^{[\bar{j},j]2} \right)$
228	$\mu^{[j]} = \max_{j \in \bar{p}^{[j]}} \mu^{[\bar{j},j]}$	$\phi^{[-1]}(\mathbf{x}) = \mathbf{x}, \Psi^{[-1]}(\Theta) = \mathbf{I}_n, \mathbf{g}^{[-1]} = \mathbf{1}_n$ $\phi_+^{[-1]2} = 0, \phi_-^{[-1]2} = \psi^{[-1]2} = \psi_-^{[-1]2} = 1$ ( $\tilde{\phi} \in \mathbb{R}_+$ arbitrary)
229	$\mathbf{f}(\mathbf{x}; \Theta) = \langle \Psi(\Theta), \phi(\mathbf{x}) \rangle_{\mathbf{g}}$	$\phi = \phi^{[D-1]}, \Psi = \Psi^{[D-1]}, \mathbf{g} = \mathbf{g}^{[D-1]}$ $\phi_+ = \phi_+^{[D-1]}, \phi_- = \phi_-^{[D-1]}, \psi = \psi^{[D-1]}, \psi_- = \psi_-^{[D-1]}$

Figure 1: Recursive definition of the global dual and bounds. See theorem 1, section 5 for details.

234 between NTK-based predictions and actual performance (Arora et al., 2019b; Lee et al., 2019). One  
 235 approach to bridging this gap is to construct higher-order or exact models. Works in this direction  
 236 include (Bai & Lee, 2019), which presented a higher-order approximation; (Bell et al., 2023), which  
 237 used a pathwise kernel; and (Shilton et al., 2023), which used an RKBS model.<sup>1</sup>

#### 239 4 HERMITE REPRESENTATION OF NEURAL ACTIVATIONS

241 The (probabilist’s) Hermite polynomials are given by (Abramowitz et al., 1972; Morse & Feshbach,  
 242 1953; Olver et al., 2010; Courant & Hilbert, 1937)  $He_k(\zeta) = (-1)^k e^{\zeta^2/2} \frac{d^k}{d\zeta^k} e^{-\zeta^2/2} \forall k \in \mathbb{N}$  and  
 243 form an orthogonal basis of  $L^2(\mathbb{R}, e^{-x^2})$ . Both models we present here make use of the Hermite  
 244 transform of the activations. For all edges  $(\bar{j}, j)$  in the network we define centered activations:

$$245 \quad \bar{\tau}^{[\bar{j},j]}(\zeta; \xi) = \tau^{[\bar{j},j]}(\xi + \zeta) - \tau^{[\bar{j},j]}(\xi)$$

247 which is simply a shifted form of the original activation  $\tau^{[\bar{j},j]}$ . By assumption  $\tau^{[\bar{j},j]} \in L^2(\mathbb{R}, e^{-\zeta^2})$ ,  
 248 so  $\bar{\tau}^{[\bar{j},j]}(\cdot, \xi) \in L^2(\mathbb{R}, e^{-\zeta^2})$  and hence the Hermite transform exists and is denoted:

$$249 \quad \begin{aligned} \bar{\tau}^{[\bar{j},j]}(\zeta; \xi) &= \sum_{k \geq 0} a_{(\xi)k}^{[\bar{j},j]} He_k(\zeta) \\ &= \sum_{k \geq 1} a_{(\xi)k}^{[\bar{j},j]} \sum_{l=1}^k \binom{k}{l} He_{k-l} \zeta^l \end{aligned} \quad (7)$$

253 where:

$$254 \quad a_{(\xi)k}^{[\bar{j},j]} = \frac{1}{\sqrt{2\pi k!}} \int_{-\infty}^{\infty} \bar{\tau}^{[\bar{j},j]}(\zeta; \xi) He_k(\zeta) e^{-\zeta^2/2} d\zeta \quad (8)$$

255 and  $He_k = He_k(0)$  are the (probabilist’s) Hermite numbers. Note that in the second form of  $\bar{\tau}^{[\bar{j},j]}$  we  
 256 use the additivity properties of the Hermite polynomials (see Appendix A).

258 For linear activations  $\tau(\zeta) = \zeta$  then trivially  $a_{(\xi)k} = \delta_{k,0}$ . For ReLU activations  $\tau(\zeta) = [\zeta]_+$ , and as  
 259 shown in Appendix A.3 (the general case  $\forall \xi \in \mathbb{R}$  is more complex - see Appendix A.3):

$$260 \quad a_{(0)k} = \begin{cases} \frac{(-1)^{p+1}}{\sqrt{2\pi}(2p-1)2^p p!} & \text{if } k = 2p, p \in \mathbb{Z}_+ \\ \frac{1}{2} \delta_{k,1} & \text{otherwise} \end{cases} \quad (9)$$

#### 264 5 GLOBAL DUAL MODEL IN REPRODUCING KERNEL BANACH SPACE

266 In this section we derive a dual model for the network described, where by global we mean not  
 267 constructed about some weight initialization. Our derivation is similar to (Shilton et al., 2023), but

268 <sup>1</sup>In a similar vein, (Bartolucci et al., 2023; Sanders, 2020; Parhi & Nowak, 2021; Unser, 2021; 2019) explore  
 269 the link to RKBS theory, though excepting (Unser, 2019) they only consider shallow networks.

based on a Hermite polynomial expansion rather than a Taylor series, making it applicable to a wider range of activation functions with fewer caveats. Our key result for this section is:<sup>2</sup>

**Theorem 1.** Let  $\mathbf{f} : \mathbb{X} \times \mathbb{W} \rightarrow \mathbb{R}^m$  be a neural network (1) satisfying our assumptions,  $\tilde{\phi} \in \mathbb{R}_+$ . Then:

$$\mathbf{f}(\mathbf{x}; \Theta) = \langle \Psi(\Theta), \phi(\mathbf{x}) \rangle_{\mathbf{g}} \quad (10)$$

with feature maps  $\Psi : \mathbb{W} \rightarrow \mathcal{W} = \overline{\text{span}}(\Psi(\mathbb{W}))$  and  $\phi : \mathbb{X} \rightarrow \mathcal{X} = \overline{\text{span}}(\phi(\mathbb{X}))$  and metric  $\mathbf{g}$  defined in Figure 1, where  $\|\Psi(\Theta)\|_2 \leq \psi$  and  $\phi_{\downarrow} \leq \|\phi(\mathbf{x})\|_2 \leq \phi = 1 \forall \Theta \in \mathbb{W}, \mathbf{x} \in \mathbb{X}$ . Moreover:

$$\|\mathbf{f}(\mathbf{x}; \Theta)\|_2 \leq \|\Psi(\Theta)\|_{\text{He}[\tau]} \|\phi(\mathbf{x})\|_2, \quad \|\Psi\|_{\text{He}[\tau]}^2 = \sup_{\mathbf{x} \in \mathbb{X}} \frac{\|\langle \Psi, \phi(\mathbf{x}) \rangle_{\mathbf{g}}\|_2^2}{\|\phi(\mathbf{x})\|_2^2} \quad (11)$$

where  $\|\Psi(\Theta)\|_{\text{He}[\tau]} \leq \psi \forall \Theta \in \mathbb{W}$ , as per definitions in Figure 1.

See Appendix B for a proof of this theorem. Intuitively, this result may be derived recursively, starting from the input node and progressing to the output, using the Hermite expansion of the activation for the edges. The operator-norm based bound (11) is included here due to the fact that the indefinite metric prevents us from naively bound  $\|\mathbf{f}(\mathbf{x}; \Theta)\|_2$  in terms of  $\phi\psi$  using the Cauchy-Schwarz inequality (as may be required e.g. when bounding Rademacher complexity).

Note that the norm-bound in Theorem 1 is defined in terms of the magnitude functions for each edge:

$$\bar{s}^{[\bar{j}, j]}(\zeta) = \sum_{k \geq 0} \left| a_{(0)k}^{[\bar{j}, j]} \right| (1 + \zeta)^k - \sum_{k \geq 0} \left| a_{(0)k}^{[\bar{j}, j]} \right| \quad (12)$$

The magnitude functions converge everywhere, are origin-crossing, monotonically increasing and superadditive on  $\mathbb{R}^+$  - for details see Appendix A.2. For linear activations  $\bar{s}(\zeta) = \zeta$ , and for ReLU activations, as shown in Appendix A.3:

$$\bar{s}(\zeta) = \frac{1}{2}\zeta \left( \text{erfi}\left(\frac{1+\zeta}{\sqrt{2}}\right) + 1 \right) + \frac{1}{\sqrt{2\pi}} \left( e^{\frac{1}{2}} - e^{\frac{1}{2}(1+\zeta)^2} \right) + \frac{1}{2} \left( \text{erfi}\left(\frac{1+\zeta}{\sqrt{2}}\right) - \text{erfi}\left(\frac{1}{\sqrt{2}}\right) \right) \quad (13)$$

## 5.1 IMPLICATION: NEURAL NETWORKS IN REPRODUCING KERNEL BANACH SPACE

A reproducing kernel Banach space is defined as:

**Definition 1** (Reproducing kernel Banach space (RKBS)). A reproducing kernel Banach space on a set  $\mathbb{X}$  is a Banach space  $\mathcal{F}$  of functions  $\mathbf{f} : \mathbb{X} \rightarrow \mathbb{Y}$  for which the point evaluation functionals  $\delta_{\mathbf{x}}(\mathbf{f}) = \mathbf{f}(\mathbf{x})$  on  $\mathcal{F}$  are continuous (i.e.  $\forall \mathbf{x} \in \mathbb{X} \exists C_{\mathbf{x}} \in \mathbb{R}_+$  such that  $\|\delta_{\mathbf{x}}(\mathbf{f})\|_2 \leq C_{\mathbf{x}} \|\mathbf{f}\|_{\mathcal{F}} \forall \mathbf{f} \in \mathcal{F}$ ).

This definition is somewhat generic, so (Lin et al., 2022)<sup>3</sup> study the special case:

$$\mathcal{B} = \{ \mathbf{f}(\cdot; \Theta) = \langle \Psi(\Theta), \Phi(\cdot) \rangle_{\mathcal{W} \times \mathcal{X}} \mid \Theta \in \mathbb{W} \} \quad (14)$$

where  $\Phi : \mathbb{X} \rightarrow \mathcal{X}$  is a data feature map,  $\Psi : \mathbb{W} \rightarrow \mathcal{W}$  is a weight feature map,  $\mathcal{X}$  and  $\mathcal{W}$  are Banach spaces, and  $\langle \cdot, \cdot \rangle_{\mathcal{W} \times \mathcal{X}} : \mathcal{W} \times \mathcal{X} \rightarrow \mathbb{R}^m$  is continuous. The following result follows from theorem 1.<sup>4</sup>

**Theorem 2.** The set  $\mathcal{F} = \{ \mathbf{f}(\cdot; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R}^m \mid \Theta \in \mathbb{W} \}$  of networks (1) satisfying our assumptions forms a RKBS of form (14), where  $\|\mathbf{f}(\cdot; \Theta)\|_{\mathcal{F}} = \|\Psi(\Theta)\|_{\text{He}[\tau]} \leq \psi$  and  $C_{\mathbf{x}} = \|\phi(\mathbf{x})\|_2 \leq \phi$ .

## 5.2 APPLICATION: BOUNDING RADEMACHER COMPLEXITY

The global dual formulation may be used to bound Rademacher complexity, which in turn bounds the uniform convergence properties of the network class (Bartlett & Mendelson, 2002; Steinwart & Christman, 2008) (that is, the rate at which the empirical risk approaches the actual risk as a function of dataset size  $N$ ). Assuming  $\mathbf{x} \sim \nu$ , the Rademacher complexity is defined as  $\mathcal{R}_N(\mathcal{F}) = \mathbb{E}_{\nu} \mathbb{E}_{\epsilon} [\sup_{f \in \mathcal{F}} \frac{1}{N} \sum_{i \in \mathbb{N}_N} \epsilon_i f(\mathbf{x}_i)]$  for Rademacher random variables  $\epsilon_i \in \{\pm 1\}$ . We have:

<sup>2</sup>When reading the norm-bounds on the feature maps in this theorem it is important to recall that  $\beta^{[j]}$  and  $\mu^{[\bar{j}, j]}$  are convenience factors representing an upper bounds on the value of  $\|\mathbf{b}^{[j]}\|_2$  and  $\|\mathbf{W}^{[\bar{j}, j]}\|_2$ , respectively. We assume these are finite, but in general their value will depend on weight-initialization, dataset complexity and regularization (if any is used).

<sup>3</sup>See e.g. (Der & Lee, 2007; Lin et al., 2022; Zhang et al., 2009; Zhang & Zhang, 2012; Song et al., 2013; Sriperumbudur et al., 2011; Xu & Ye, 2014) for other perspectives.

<sup>4</sup>Note that the RKBS defined in theorem 2 is non-reflexive, which appears to rule out a trivial representor theory based on this dual in the global case (Lin et al., 2022).

**Theorem 3.** *The set  $\mathcal{F} = \{f(\cdot; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R} \mid \Theta \in \mathbb{W}\}$  of networks (1) satisfying our assumptions has Rademacher complexity bounded by  $\mathcal{R}_N(\mathcal{F}) \leq \frac{1}{\sqrt{N}} \phi \underline{\psi} = \frac{1}{\sqrt{N}} \underline{\psi}$  (definitions as per Figure 1).*

The proof follows the usual template (see e.g. (Bartlett & Mendelson, 2002)) using our feature map, with (11) used instead of the Cauchy-Schwarz inequality (Appendix D). Intuitively, we may think of the recursive bounds  $\underline{\psi}^{[j]}$ ,  $\underline{\psi}^{[\tilde{j}, j]}$  on the weight feature map in terms of signal flow in a electrical circuit, precisely (with reference to figure 1):

1. A signal  $\underline{\psi}^{[-1]} = 1$  enters the network at the input node  $\tilde{j} = -1$ .
2. The outgoing edge from node  $(\tilde{j}, j)$  amplifies this signal:

$$\underline{\psi}^{[\tilde{j}, j]} = \sqrt{\sup_{\substack{\phi^{[\tilde{j}]} \leq \phi^{[j]} \leq 1 \\ -\underline{\psi}^{[\tilde{j}]} \leq \underline{\psi}^{[j]} \leq \underline{\psi}^{[\tilde{j}]}}} \left\{ \frac{\tilde{\phi}^2 \tau^{[\tilde{j}, j]} (\phi^{[j]} \underline{\psi}^{[\tilde{j}]})^2}{\tilde{s}^{[\tilde{j}, j]} (\tilde{s}^{[\tilde{j}, j]} - 1) (\tilde{\phi}^2) \phi^{[j]2}} \right\}} \quad (15)$$

3. Subsequent nodes  $j$  combine signals from incoming edges  $(\tilde{j}, j)$  into an offset weighted sum:

$$\underline{\psi}^{[j]2} = (\gamma^2 + 1) \left( \left( \beta^{[j]} + \frac{1}{\gamma} \mu^{[j]} |\tau^{[\tilde{j}, j]}(0)| \right)^2 + \sum_{\tilde{j} \in \mathbb{P}^{[j]}} \frac{\tilde{H}^{[j]}}{H^{[j]}} \underline{\psi}^{[\tilde{j}, j]2} \mu^{[\tilde{j}, j]2} \right) \quad (16)$$

4. The signal propagates (steps 2-3) to the output  $D - 1$ . The overall output is  $\underline{\psi}^{[D-1]}$ .

For Lipschitz activations (ie. most activations) we can simplify step 2 with the following theorem:

**Theorem 4.** *For  $L$ -Lipschitz neural activations, in the limit  $\tilde{\psi} \rightarrow 0_+$  (recall that  $\tilde{\psi} \in \mathbb{R}_+$ ), we have that  $\underline{\psi}^{[j]2} = (\gamma^2 + 1) \left( (\beta^{[j]} + \frac{1}{\gamma} \mu^{[j]} |\tau^{[\tilde{j}, j]}(0)| \right)^2 + L^2 \sum_{\tilde{j} \in \mathbb{P}^{[j]}} \frac{\tilde{H}^{[j]}}{H^{[j]}} \underline{\psi}^{[\tilde{j}, j]2} \mu^{[\tilde{j}, j]2} \forall j$ , which in the unbiased case this simplifies to  $\underline{\psi}^{[j]2} = L^2 \sum_{\tilde{j} \in \mathbb{P}^{[j]}} \frac{\tilde{H}^{[j]}}{H^{[j]}} \underline{\psi}^{[\tilde{j}, j]2} \mu^{[\tilde{j}, j]2}$ .*

from which we obtain the corollary for unbiased networks:<sup>5</sup>

**Corollary 5.** *Let  $\mathcal{F} = \{f(\cdot; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R} \mid \Theta \in \mathbb{W}\}$  be the set of networks (1) with zero bias  $\gamma = \tau^{[\tilde{j}, j]}(0) = 0$ , and  $\mu^{[\tilde{j}, j]2} \leq \frac{H^{[j]}}{L^2 \tilde{H}^{[j]}} \forall j, \tilde{j} \in \mathbb{P}^{[j]}$ , Rademacher complexity is bounded  $\mathcal{R}_N(\mathcal{F}) \leq \frac{1}{\sqrt{N}}$ .*

This shows that Rademacher complexity is depth-independent for sufficiently small network weights, but exponential in depth (the longest path in the network) for large weights; and that width dependence will scale with the product of  $\mu^{[\tilde{j}, j]}$  along the longest path. To gain further insight it is worth considering the Rademacher complexity of unbiased, randomly networks  $W_{i_j, i_j}^{[j]} \sim \mathcal{N}(0, \sigma^{[j]2})$ . We will consider the depth and width dependence of the complexity bound separately:

- **Depth dependence:** from (4) and corollary 5, the Rademacher complexity will be depth-independent and satisfy  $\mathcal{R}_N(\mathcal{F}) \leq \frac{1}{\sqrt{N}}$  whp  $\geq 1 - \epsilon$  if:

$$\sigma^{[j]2} \leq \frac{1}{L^2 \left( \tilde{H}^{[j]} + 2 \frac{\tilde{H}^{[j]}}{\sqrt{H^{[j]}}} \ln \left( \frac{D H^{[j]}}{2\epsilon} \right) + 2 \frac{\tilde{H}^{[j]}}{H^{[j]}} \ln \left( \frac{D H^{[j]}}{\epsilon} \right) \right)} \quad (17)$$

Note that this is a modified He initialization accounting for neural activation slope (through  $L$ ) and correction terms for network topology  $\frac{\tilde{H}^{[j]}}{H^{[j]}}$ , node count  $D$  and fan-out  $H^{[j]}$ . A similar modified Glorot initialization follows trivially.

- **Width dependence:** ignoring depth, we observe that for Glorot (and modified He/Glorot) initialization  $\mathcal{R}_N(\mathcal{F}) \sim \mathcal{O}(1)$ . LeCun and He initialization behave similarly if  $H^{[j]} \asymp \tilde{H}^{[j]}$ , but LeCun initialization may scale arbitrarily at the output node ( $H^{[D-1]} = m$ , while  $\tilde{H}^{[j]}$  is arbitrary), and He initialization may be analogously badly behaved at the input.

With regard to our ReLU and ResNet examples (see section 2), both will have Rademacher complexity  $\mathcal{R}_N(\mathcal{F}) \leq \frac{1}{\sqrt{N}}$  if the spectral norm of all weight matrices is less than 1 (ReLU) or  $\frac{1}{2}$  (ResNet). This will also hold for random ReLU/ResNet networks whp for modified He (17) (or Glorot) initialization.

<sup>5</sup>In general the condition in the corollary is  $(\gamma^2 + 1) \left( (\beta^{[j]} + \frac{1}{\gamma} \mu^{[j]} |\tau^{[\tilde{j}, j]}(0)| \right)^2 + L^2 \sum_{\tilde{j} \in \mathbb{P}^{[j]}} \frac{\tilde{H}^{[j]}}{H^{[j]}} \underline{\psi}^{[\tilde{j}, j]2} \mu^{[\tilde{j}, j]2} \leq 1$ .

	Incoming Edge Feature Map	Node Feature Map
378		
379		
380	$\tilde{\phi}_{\Delta}^{[j]}(\mathbf{x}) = \frac{1}{T_{(\tilde{\omega})}^{[j]}} \left[ \frac{1}{\eta^k} \left[ \left( \sqrt{\frac{[\tilde{\omega}]^2}{(\tilde{\omega})\eta}} \phi_{\Delta}^{[j]}(\mathbf{x}) \right)^{\otimes l} \right]_{1 \leq l \leq k} \right]_{k \geq 1}$	$\phi_{\Delta}^{[j]}(\mathbf{x}) = \sqrt{\frac{1}{\gamma^2 + 1}} \begin{bmatrix} \frac{\gamma}{\sqrt{2}} \left[ \sqrt{\frac{H^{[j]}}{H^{[j]}}} \frac{1}{\tilde{\omega}^{[j]}} \tilde{\mathbf{x}}^{[j]} \right]_{j \in \tilde{\mathbb{P}}^{[j]}} \\ \frac{1}{\sqrt{2}} \left[ \sqrt{\frac{H^{[j]}}{H^{[j]}}} \tilde{\phi}_{\Delta}^{[j]}(\mathbf{x}) \right]_{j \in \tilde{\mathbb{P}}^{[j]}} \end{bmatrix}$
381		
382	$\tilde{\Psi}_{\Delta}^{[j]}(\Theta) = T_{(\tilde{\omega})}^{[j]} \left[ \eta^k \left[ \left( \sqrt{\frac{1}{\rho^{[j]}} \Psi_{\Delta}^{[j]}(\Theta)} \right)^{\otimes l} \right]_{1 \leq l \leq k} \right]_{k \geq 1}$	$\Psi_{\Delta}^{[j]}(\Theta) = \sqrt{\gamma^2 + 1} \begin{bmatrix} \sqrt{2} \text{diag}_{j \in \tilde{\mathbb{P}}^{[j]}} \left( \sqrt{\frac{H^{[j]}}{H^{[j]}}} \tilde{\omega}^{[j]} \Delta \mathbf{W}^{[j]} \right) \\ \sqrt{2} \text{diag}_{j \in \tilde{\mathbb{P}}^{[j]}} \left( \sqrt{\frac{H^{[j]}}{H^{[j]}}} \tilde{\Psi}_{\Delta}^{[j]}(\Theta) (\mathbf{W}^{[j]} + \Delta \mathbf{W}^{[j]}) \right) \end{bmatrix}$
383		
384	$\tilde{\mathbf{G}}_{\Delta}^{[j]}(\mathbf{x}) = \begin{bmatrix} a_{(\tilde{\omega})}^{[j]} \\ \left( \frac{k}{i} \right) \text{He}_{k-1} \mathbf{G}_{\Delta}^{[j]}(\mathbf{x})^{\otimes l} \end{bmatrix}_{1 \leq l \leq k} \Big _{k \geq 1}$	$\mathbf{G}_{\Delta}^{[j]}(\mathbf{x}) = \begin{bmatrix} \mathbf{1}_{H^{[j]}}^T \\ \text{diag}_{j \in \tilde{\mathbb{P}}^{[j]}} \left( \mathbf{1}_{H^{[j]}} \mathbf{1}_{H^{[j]}}^T \right) \\ \text{diag}_{j \in \tilde{\mathbb{P}}^{[j]}} \left( \tilde{\mathbf{G}}_{\Delta}^{[j]}(\mathbf{x}) \mathbf{1}_{H^{[j]}} \mathbf{1}_{H^{[j]}}^T \right) \end{bmatrix}$
385		
386		
387	<b>Bounds</b> $\left\  \tilde{\phi}_{\Delta}^{[j]}(\mathbf{x}) \odot \tilde{\mathbf{G}}_{\Delta}^{[j]}(\mathbf{x}) \right\ _2 \leq 1 \quad \forall i_j$	$\left\  \phi_{\Delta}^{[j]}(\mathbf{x}) \odot \mathbf{G}_{\Delta}^{[j]}(\mathbf{x}) \right\ _2 \leq \phi_{\Delta}^{[j]2} = 1 \quad \forall i_j$
388	$\left\  \tilde{\Psi}_{\Delta}^{[j]}(\Theta) \right\ _2 \leq \tilde{\psi}_{\Delta}^{[j]2} = T_{(\tilde{\omega})}^{[j]2} s_{\eta} \left( \frac{\gamma^2 + 1}{\rho^{[j]}} \psi_{\Delta}^{[j]2} \right)$	$\left\  \Psi_{\Delta}^{[j]}(\Theta) \right\ _2 \leq \psi_{\Delta}^{[j]2} = (\gamma^2 + 1) \left( \beta^{[j]2} + 2 \sum_{j \in \tilde{\mathbb{P}}^{[j]}} \frac{H^{[j]}}{H^{[j]}} \left( \tilde{\omega}^{[j]2} \mu_{\Delta}^{[j]2} + \left( \mu^{[j]2} + \mu_{\Delta}^{[j]2} \right) \tilde{\psi}_{\Delta}^{[j]2} \right) \right)$
389	$\tilde{\omega}^{[j]} = \max_{j \in \tilde{\mathbb{P}}^{[j]}} \tilde{\omega}^{[j]}, T_{(\tilde{\omega})}^{[j]} = \max_{j \in \tilde{\mathbb{P}}^{[j]}} T_{(\tilde{\omega})}^{[j]}, \tilde{\psi}_{\Delta}^{[j]} = \max_{j \in \tilde{\mathbb{P}}^{[j]}} \tilde{\psi}_{\Delta}^{[j]}$	$\phi_{\Delta}^{[-1]}(\mathbf{x}) = \mathbf{0}, \Psi_{\Delta}^{[-1]}(\Theta) = \mathbf{G}_{\Delta}^{[-1]}(\mathbf{x}) = \mathbf{1}_{0 \times n}, \phi_{\Delta}^{[-1]2} = \psi_{\Delta}^{[-1]2} = 0$
390	$\mathbf{f}(\mathbf{x}; \Theta + \Delta\Theta) = \mathbf{f}(\mathbf{x}; \Theta) + \Delta\mathbf{f}(\mathbf{x}; \Delta\Theta)$	$\phi_{\Delta} = \phi_{\Delta}^{[D-1]}, \Psi_{\Delta} = \Psi_{\Delta}^{[D-1]}, \mathbf{G}_{\Delta} = \mathbf{G}_{\Delta}^{[D-1]}$
391	$\Delta\mathbf{f}(\mathbf{x}; \Delta\Theta) = \langle \Psi_{\Delta}(\Delta\Theta), \phi_{\Delta}(\mathbf{x}) \rangle_{\mathbf{G}_{\Delta}(\mathbf{x})}$	$\phi_{\Delta} = \phi_{\Delta}^{[D-1]}, \psi_{\Delta} = \psi_{\Delta}^{[D-1]}$

Figure 2: Recursive definition of local dual and bounds. See theorem 6, section 6 for details.

More generally in uniform convergence analysis we must consider how the weight-norm  $\mu^{[\tilde{j}, j]}$  evolves or increases during training. It is difficult to draw firm conclusions about this without delving into the specifics of training, however in the lazy regime, or otherwise given sufficiently strong regularization, we would expect that this norm-bound should remain close to its initialization value, potentially indicating good uniform-convergence behavior for a wide class of neural networks.

## 6 LOCAL DUAL MODEL IN REPRODUCING KERNEL HILBERT SPACE

When considering training or network adaptation it is better to model the change in the network rather than the network in-toto. To this end, in this section we present an exact (non-approximate) local RKHS model. Let  $\Theta$  be the initial weight and biases and  $\Delta\Theta$  the change in weights and biases - so  $\Delta\Theta$  might be a training step, multiple steps, or even the complete training process after random initialization. Let  $\tilde{\mathbf{x}}^{[j]}, \mathbf{x}^{[j]}$  denote the pre-activation (input) and post-activation (output) of node  $j$  with initial weights  $\Theta$  given input  $\mathbf{x}$ ; and  $\Delta\tilde{\mathbf{x}}^{[j]}, \Delta\mathbf{x}^{[j]}$  the change in same due to the change in weights  $\Delta\Theta$ . The change in network operation is denoted  $\Delta\mathbf{f} : \mathbb{X} \times \mathbb{W}_{\Delta} \rightarrow \mathbb{R}^m$ :

$$\mathbf{f}(\mathbf{x}; \Theta + \Delta\Theta) = \mathbf{f}(\mathbf{x}; \Theta) + \Delta\mathbf{f}(\mathbf{x}; \Delta\Theta), \quad (18)$$

For the purposes of this analysis we augment our previous assumptions with in section 2 with:

4. **Bounded activation:**  $\|\tilde{\mathbf{x}}^{[\tilde{j}, j]}\|_2 \leq \tilde{\omega}^{[\tilde{j}, j]} \forall \tilde{j} \in \tilde{\mathbb{P}}^{[j]}$  (note that  $\tilde{\omega}^{[\tilde{j}, j]} \leq \tilde{\phi}^{[\tilde{j}, j]} \tilde{\psi}^{[\tilde{j}, j]}$ ).
5. **Bounded weight and bias steps:**  $\|\Delta\mathbf{W}^{[\tilde{j}, j]}\|_2 \leq \mu_{\Delta}^{[\tilde{j}, j]} \leq \mu^{[\tilde{j}, j]}, \|\Delta\mathbf{b}^{[j]}\|_2 \leq \beta_{\Delta}^{[j]}$ .

satisfying (20) (details in and after Theorem 6). Note that, unlike the parameters  $\mu^{[\tilde{j}, j]}, \beta^{[j]}$  in our global model which may be arbitrarily large and as such do not place any restriction on the network which may be modeled, the parameters  $\mu_{\Delta}^{[\tilde{j}, j]}, \beta_{\Delta}^{[j]}$  bounding step must satisfy (20) and are constrained themselves and subsequently constrain the size of step that can be modeled by the local model.

We define  $\mathbb{W}_{\Delta}$  to be the set of all weight-steps satisfying these assumptions. The parameter  $\tilde{\omega}^{[\tilde{j}, j]}$  is a bound on the magnitude of the output of edge ( $\tilde{j} \rightarrow j$ ) in our initial network  $\forall \mathbf{x} \in \mathbb{X}$ . With this prequel, we have the following local analogue of theorem 1 (see proof in Appendix C):<sup>6</sup>

**Theorem 6.** Let  $\Delta\mathbf{f} : \mathbb{X} \times \mathbb{W}_{\Delta} \rightarrow \mathbb{R}^m$  be the change in neural network operation (18). Then:

$$\Delta\mathbf{f}(\mathbf{x}; \Delta\Theta) = \langle \Psi_{\Delta}(\Delta\Theta), \phi_{\Delta}(\mathbf{x}) \rangle_{\mathbf{G}_{\Delta}(\mathbf{x})} \quad (19)$$

with feature maps  $\phi_{\Delta} : \mathbb{X} \rightarrow \mathcal{X}_{\Delta} = \overline{\text{span}}(\phi_{\Delta}(\mathbb{X}))$ ,  $\Psi_{\Delta} : \mathbb{W}_{\Delta} \rightarrow \mathcal{W}_{\Delta} = \overline{\text{span}}(\Psi_{\Delta}(\mathbb{W}_{\Delta}))$  and metric  $\mathbf{G}_{\Delta}(\mathbf{x})$  as per Figure 2, where  $\|\phi_{\Delta}(\mathbf{x}) \odot \mathbf{G}_{\Delta}^{i_{D-1}}(\mathbf{x})\|_2 \leq \phi_{\Delta} \forall i_{D-1}$  and  $\|\Psi_{\Delta}(\Delta\Theta)\|_2 \leq$

<sup>6</sup>The decision to use a position dependent metric  $\mathbf{G}_{\Delta}$  here is largely stylistic. We could of course absorb  $\mathbf{G}_{\Delta}$  into  $\phi_{\Delta}$  without substantively changing our results.

$\psi_\Delta \forall \mathbf{x} \in \mathbb{X}, \Delta\Theta \in \mathbb{W}_\Delta$ . Moreover  $\|\Psi_\Delta(\Delta\Theta)\|_2 \leq \psi_\Delta = S_\eta^2 < 1$  if  $\forall j$ :

$$\mu_\Delta^{[\tilde{j},j]^2} + \frac{1}{2\bar{\rho}^{[\tilde{j}]\omega^{[\tilde{j}]}]}} \beta_\Delta^{[j]^2} \leq \frac{u^{[j]^2}}{4\bar{\rho}^{[\tilde{j}]\omega^{[\tilde{j}]}]}} : u^{[j]^2} = \min_{\tilde{j}:j \in \mathbb{P}^{[\tilde{j}]}} \rho_{(\bar{\omega}^{[\tilde{j}]\omega^{[\tilde{j}]}})_\eta}^{[\tilde{j},\tilde{j}]^2} \left\{ R_\eta^2, \hat{s}_\eta^{-1} \left( \frac{u^{[\tilde{j}]^2}}{8\bar{\rho}^{[\tilde{j}]\mu^{[\tilde{j}]^2}}} \right) \right\} \quad (20)$$

We emphasise that this exactly models the change in the neural network without approximation so long as conditions in the theorem are met. This is in contrast to the NTK model, which is a first-order approximation whose accuracy will decrease as the step-size increases (e.g. for narrower networks).

In constructing Theorem 6 we use the *rectified activation functions* and their envelopes, respectively:

$$\begin{aligned} \hat{\tau}_\eta^{[\tilde{j},j]}(\zeta; \xi, \xi') &= \sum_{k \geq 1} \frac{a_{(\xi)k}^{[\tilde{j},j]} a_{(\xi')k}^{[\tilde{j},j]}}{\eta^{2k}} \sum_{l=1}^k \binom{k}{l}^2 \text{He}_{k-l}^2 \zeta^l \\ \hat{\tau}_\eta^{[\tilde{j},j]}(\zeta; \omega^{[\tilde{j},j]}) &= \sup_{|\xi|, |\xi'| \leq \omega^{[\tilde{j},j]}} \hat{\tau}_\eta^{[\tilde{j},j]}(\zeta; \xi, \xi') = \sup_{|\xi| \leq \omega^{[\tilde{j},j]}} \hat{\tau}_\eta^{[\tilde{j},j]}(\zeta; \xi, \xi) \end{aligned} \quad (21)$$

where  $\xi, \xi' \in [-\omega^{[\tilde{j},j]}, \omega^{[\tilde{j},j]}]$  are the centers of the rectified activation functions (the initial activation for some input about which our model is constructed) and  $\eta \in (0, 1)$  is fixed. Unlike the magnitude functions, the rectified activations have a finite ROC  $|\hat{\tau}_\eta^{[\tilde{j},j]}(\zeta; \xi, \xi')| \leq T_{(\xi, \xi')_\eta}^{[\tilde{j},j]^2} \forall |\zeta| \leq \rho_{(\xi, \xi')_\eta}^{[\tilde{j},j]^2}$  and likewise  $|\hat{\tau}_\eta^{[\tilde{j},j]}(\zeta; \omega^{[\tilde{j},j]})| \leq T_{(\omega^{[\tilde{j},j]})_\eta}^{[\tilde{j},j]^2} \forall |\zeta| \leq \rho_{(\omega^{[\tilde{j},j]})_\eta}^{[\tilde{j},j]^2}$  (see Appendix A.2). The rectified activation envelopes are origin crossing, monotonically increasing and superadditive. We also define:

$$\hat{s}_\eta(\zeta) = \sum_{k \geq 1} \eta^{2k} \sum_{1 \leq l \leq k} \zeta^l = \frac{\zeta}{1-\zeta} \left( \frac{\eta^2}{1-\eta^2} - \frac{\zeta \eta^2}{1-\zeta \eta^2} \right)$$

which converges as given  $\forall |\zeta| < R_\eta^2 < 1$ , whereon  $|\hat{s}_\eta(\zeta)| \leq S_\eta^2 = \frac{R_\eta^2}{1-R_\eta^2} \left( \frac{\eta^2}{1-\eta^2} - \frac{\eta^2 R_\eta^2}{1-\eta^2 R_\eta^2} \right)$ .

It is difficult to obtain a closed-form expression for the rectified activation or its envelope for the ReLU, but they are relatively straightforward to calculate, as are their convergence bounds. Figure 3 in Appendix A.3 shows a sample of various rectified activations for the ReLU with different centers.

## 6.1 IMPLICATION: NEURAL NETWORK CHANGE IN REPRODUCING KERNEL HILBERT SPACE

A vector-valued (v-v) reproducing kernel Hilbert space is defined as follows (Aronszajn, 1950; Steinwart & Christman, 2008; Shawe-Taylor & Cristianini, 2004; Mercer, 1909; Micchelli & Pontil, 2005; Caponnetto & De Vito, 2007; Reiser & Burkhart, 2007; Carmeli et al., 2005; Schwartz, 1964):

**Definition 2** (Reproducing kernel Hilbert space (RKHS)). A v-v reproducing kernel Hilbert space  $\mathcal{H}$  on a set  $\mathbb{X}$  is a Hilbert space  $\mathcal{F}$  of functions  $\mathbf{f} : \mathbb{X} \rightarrow \mathbb{R}^m$  for which the point evaluation functionals  $\delta_{\mathbf{x}}(\mathbf{f}) = \mathbf{f}(\mathbf{x})$  on  $\mathcal{F}$  are continuous ( $\forall \mathbf{x} \in \mathbb{X} \exists C_{\mathbf{x}} \in \mathbb{R}_+$  s.t.  $\|\delta_{\mathbf{x}}(\mathbf{f})\|_2 \leq C_{\mathbf{x}} \|\mathbf{f}\|_{\mathcal{F}} \forall \mathbf{f} \in \mathcal{F}$ ).

For an RKHS, Reisz representer theory implies that  $\forall \mathbf{x} \in \mathbb{X} \exists$  unique  $\mathbf{K}_{\mathbf{x}} \in \mathcal{F} \times \mathbb{R}^m$  such that  $\langle \mathbf{f}(\mathbf{x}), \mathbf{v} \rangle = \langle \mathbf{f}, \mathbf{K}_{\mathbf{x}} \mathbf{v} \rangle_{\mathcal{H}} \forall \mathbf{v} \in \mathbb{R}^m$ . From this, the kernel  $\mathbf{K} : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}^{m \times m}$  is defined as:

$$\mathbf{K}(\mathbf{x}, \mathbf{x}') = \left[ \left\langle \mathbf{K}_{\mathbf{x}} \delta_{(i_{D-1})}^{[D-1]}, \mathbf{K}_{\mathbf{x}'} \delta_{(i'_{D-1})}^{[D-1]} \right\rangle_{\mathcal{H}} \right]_{i_{D-1}, i'_{D-1}}, \quad \text{where } \delta_{(k)}^{[j]} = [\delta_{k, i_j}]_{i_j}$$

Moore-Aronszajn theorem allows us to run the argument in reverse: any symmetric, positive definite  $\mathbf{K} : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}^{m \times m}$  uniquely defines an RKHS,  $\mathcal{H}_{\mathbf{K}}$  for which  $\mathbf{K}$  is the kernel. From theorem 6:

**Theorem 7.** The set  $\mathcal{F}_\Delta = \{\Delta \mathbf{f}(\cdot; \Delta\Theta) : \mathbb{R}^n \rightarrow \mathbb{R}^m | \Delta\Theta \in \mathbb{W}_\Delta\}$  of changes in network behavior satisfying our assumptions, including the bound, lies in an RKHS  $\mathcal{H}_{\mathbf{K}}$  (that is,  $\mathcal{F}_\Delta \subset \mathcal{H}_{\mathbf{K}}$ ) with kernel  $\mathbf{K} = \mathbf{I}_m K_{\text{LiNK}}$ , where  $K_{\text{LiNK}} = K^{[D-1]}$ , is the Local-intrinsic Neural Kernel (LiNK),  $\forall j$ :

$$K^{[j]}(\mathbf{x}, \mathbf{x}') = \frac{\rho^{[j]}}{\gamma^{2+1}} \mathbb{E}_{\tilde{j} \in \mathbb{P}^{[j]}} \left[ \frac{H^{[\tilde{j}]}}{\bar{\omega}^{[\tilde{j},j]^2}} \Sigma^{[\tilde{j},j]}(\mathbf{x}, \mathbf{x}') + \frac{H^{[\tilde{j}]}}{T^{[\tilde{j},j]^2}} \mathbb{E}_{i_{\tilde{j}}} \left[ \hat{\tau}_\eta^{[\tilde{j},j]} \left( \rho_{(\bar{\omega})_\eta}^{[\tilde{j},j]^2} K^{[\tilde{j}]}(\mathbf{x}, \mathbf{x}'); x_{i_{\tilde{j}}}^{[\tilde{j}]}, x'_{i_{\tilde{j}}}^{[\tilde{j}]} \right) \right] \right] \quad (22)$$

and  $K^{[-1]}(\mathbf{x}, \mathbf{x}') = 0$  and  $\Sigma^{[j]}(\mathbf{x}, \mathbf{x}')$  is the NNGP kernel. Moreover:

$$\begin{aligned} \lim_{\eta \rightarrow 1} \mathbb{E}_{i_{\tilde{j}}} \left[ \hat{\tau}_\eta^{[\tilde{j},j]} \left( \rho_{(\bar{\omega})_\eta}^{[\tilde{j},j]^2} K^{[\tilde{j}]}(\mathbf{x}, \mathbf{x}'); x_{i_{\tilde{j}}}^{[\tilde{j}]}, x'_{i_{\tilde{j}}}^{[\tilde{j}]} \right) \right] &= \sum_{q \geq 1} \theta_q^{[\tilde{j},j]}(\mathbf{x}, \mathbf{x}') \left( \rho_{(\bar{\omega})_\eta}^{[\tilde{j},j]^2} K^{[\tilde{j}]}(\mathbf{x}, \mathbf{x}') \right)^q \\ \theta_q^{[\tilde{j},j]}(\mathbf{x}, \mathbf{x}') &= \mathbb{E}_{i_{\tilde{j}}} \left[ \frac{1}{q!} \tau^{[\tilde{j},j]}(q) \left( x_{i_{\tilde{j}}}^{[\tilde{j}]} \right) \frac{1}{q!} \tau^{[\tilde{j},j]}(q) \left( x'_{i_{\tilde{j}}}^{[\tilde{j}]} \right) \right] \end{aligned} \quad (23)$$

(here  $\theta_q^{[\tilde{j},j]}(\mathbf{x}, \mathbf{x}')$  is the raw covariance of the  $q^{\text{th}}$  derivative of link  $(\tilde{j}, j)$ 's activation given  $\mathbf{x}, \mathbf{x}'$ .)

The proof requires two steps - step one is to apply the kernel trick (after some preliminaries), while step two uses Mertens' theorem to obtain the final result - see Appendix C.3 for details. We observe that the NTK is essentially (with some additional scaling factors) a first-order (in  $q$ ) approximation of the LiNK - if we take the limit  $\eta \rightarrow 1$  then to first order the LiNK is approximately:

$$K^{[j]}(\mathbf{x}, \mathbf{x}') \approx \frac{p^{[j]}}{\gamma^2 + 1} \mathbb{E}_{\tilde{\gamma} \in \mathbb{P}^{[j]}} \left[ \frac{H^{[\tilde{\gamma}]}}{\tilde{\omega}^{[\tilde{\gamma}, j]^2}} \Sigma^{[\tilde{\gamma}, j]}(\mathbf{x}, \mathbf{x}') + \frac{\rho^{[\tilde{\gamma}, j]^2} H^{[\tilde{\gamma}]}}{T^{[\tilde{\gamma}, j]^2} (\tilde{\omega}_\eta)^{[\tilde{\gamma}, j]^2}} \theta_1^{[\tilde{\gamma}, j]}(\mathbf{x}, \mathbf{x}') K^{[\tilde{\gamma}]}(\mathbf{x}, \mathbf{x}') \right]$$

recalling  $K_{\text{LiNK}} = K^{[D-1]}$ , where:

$$\theta_1^{[\tilde{\gamma}, j]}(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{i_{\tilde{\gamma}}} \left[ \tau^{[\tilde{\gamma}, j](1)}(x_{i_{\tilde{\gamma}}}^{[\tilde{\gamma}]}) \tau^{[\tilde{\gamma}, j](1)}(x'_{i_{\tilde{\gamma}}}{}^{[\tilde{\gamma}]}) \right]$$

which is essentially the NTK (6) with some additional scale factors. Assuming random initialization the LiNK is well-defined for almost all  $\mathbf{x} \in \mathbb{X}$  if  $\tau^{[\tilde{\gamma}, j]} \in \mathcal{C}^\infty$  for almost all  $\mathbf{x} \in \mathbb{X}$ . Note however that  $\mathcal{F}_\Delta \subset \mathcal{H}_K$  - ie.  $\mathcal{F}_\Delta$  is not an RKHS in general, but rather a subspace inside of one. Nor can we meaningfully replace  $\mathcal{F}_\Delta$  with its span or completion, as this will contain elements that do not correspond to physically realizable networks. This is clear from Figure 2, where the weight-feature map  $\Psi_\Delta$  maps the network weights onto a (non-flat) subspace of  $\ell^2(\mathbb{N})^m$ , no column of which coincides with the subspace of same onto which  $\phi_\Delta$  maps input space. Thus in general the LiNK cannot be naively used as a basis for a representer theory in terms of the training dataset.

## 6.2 APPLICATION: BOUNDING RADEMACHER COMPLEXITY FOR ADAPTATION

Like the global model, an obvious application of the local dual model is the bounding of Rademacher complexity. The following result may be viewed as the local analogue of our previous bound:

**Theorem 8.** *The set  $\mathcal{F}_\Delta = \{\Delta f(\cdot; \Delta\Theta) : \mathbb{R}^n \rightarrow \mathbb{R} \mid \Delta\Theta \in \mathbb{W}_\Delta\}$  of change in neural-network operation satisfying (20) has Rademacher complexity  $\mathcal{R}_N(\mathcal{F}) \leq \frac{1}{\sqrt{N}} \phi_\Delta \psi_\Delta$  (defined in Figure 2).*

The proof follows the template of (Bartlett & Mendelson, 2002) using the local feature map and the Cauchy-Schwarz inequality (see Appendix D). Assuming an unbiased network  $\phi_\Delta = 1$ , and the Rademacher complexity bound is determined by the recursive equation:

$$\psi_\Delta^{[j]^2} = (\gamma^2 + 1) \left( \beta_\Delta^{[j]^2} + 2 \sum_{\tilde{\gamma} \in \mathbb{P}^{[j]}} \frac{\tilde{H}^{[\tilde{\gamma}]}}{H^{[\tilde{\gamma}]}} \left( \tilde{\omega}^{[\tilde{\gamma}, j]^2} \mu_\Delta^{[\tilde{\gamma}, j]^2} + \left( \mu^{[\tilde{\gamma}, j]^2} + \mu_\Delta^{[\tilde{\gamma}, j]^2} \right) T_{(\tilde{\omega}_\eta)^{[\tilde{\gamma}, j]^2}} \hat{s}_\eta \left( \frac{\gamma^2 + 1}{\rho^{[\tilde{\gamma}, j]^2}} \psi_\Delta^{[\tilde{\gamma}]^2} \right) \right) \right)$$

The width-dependence of the bound is dependent on the width-dependence of  $\mu^{[\tilde{\gamma}, j]}$  and  $\mu_\Delta^{[j]}$ , but unfortunately there appears to be an unavoidably exponential depth-dependency not present in the global model as  $\hat{s}_\eta$  is positive and increasing on  $\mathbb{R}^+$ . In future work we hope to use this theorem to explain how methods such as LoRA (Hu et al., 2021) achieve better performance in terms of uniform convergence properties with restricted weight update rank (and hence the spectral norm of the weight-matrix changes). Moreover it may be interesting in future investigation to explore if spectral analysis of the LiNK could be used to bound *local* Rademacher complexity (Cortes et al., 2013; Bartlett et al., 2005), as previous investigations in RKHS using this approach give bounds up to  $\mathcal{O}(\frac{1}{N})$ .

## 7 CONCLUSIONS AND FUTURE DIRECTIONS

In this paper we have presented two models of neural networks and neural network training for neural network of arbitrary width, depth and topology. First we presented an exact (non-approximated) RKBS model of the overall network in the form of a bilinear product between a data- and weight-feature map. We have used this model to construct a bound on Rademacher complexity, and for Lipschitz activations we have given conditions under which the Rademacher complexity is depth-independent, and how different initialization schemes can achieve  $\mathcal{R}_N(\mathcal{F}) \leq \frac{1}{\sqrt{N}}$ . The second model we have presented models the *change* in the neural network due to a bounded change in weights and biases. This model cast the change in RKHS with the local-intrinsic neural kernel (LiNK). We have shown that this can be used to bound Rademacher complexity for network adaptation. We have also discussed the role of weight initialization and implications for feedforward ReLU networks and residual networks (ResNets), and presented the local intrinsic neural kernel for the ResNet.

## REFERENCES

- 540  
541  
542 Milton Abramowitz, Irene A. Stegun, and Donald A. McQuarrie. *Handbook of Mathematical*  
543 *Functions with Formulas, Graphs, and Mathematical Tables*. Dover, 1972.
- 544 Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparameterized  
545 neural networks, going beyond two layers. *arXiv preprint arXiv:1811.04918*, 2018.
- 546  
547 Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-  
548 parameterization. In *International Conference on Machine Learning*, pp. 242–252. PMLR, 2019.
- 549 N. Aronszajn. Theory of reproducing kernels. *Transactions of the American Mathematical Society*,  
550 68:337–404, Jan–Jun 1950.
- 551  
552 Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang. Stronger generalization bounds for  
553 deep nets via a compression approach. In *Proceedings of ICML*, 2018.
- 554  
555 Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of  
556 optimization and generalization for overparameterized two-layer neural networks. In *International*  
557 *Conference on Machine Learning*, pp. 322–332. PMLR, 2019a.
- 558  
559 Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On  
560 exact computation with an infinitely wide neural net. In *Advances in Neural Information Processing*  
561 *Systems*, pp. 8139–8148, 2019b.
- 562  
563 Francis Bach. On the equivalence between kernel quadrature rules and random feature expansions.  
564 *The Journal of Machine Learning Research*, 18(1):714–751, 2017.
- 565  
566 Francis R. Bach. Breaking the curse of dimensionality with convex neural networks. *CoRR*,  
567 abs/1412.8690, 2014. URL <http://arxiv.org/abs/1412.8690>.
- 568  
569 Yu Bai and Jason D. Lee. Beyond linearization: On quadratic and higher-order approximation of  
570 wide neural networks. *arXiv preprint arXiv:1910.01619*, 2019.
- 571  
572 P. L. Bartlett and S. Mendelson. Rademacher and gaussian complexities: Risk bounds and structural  
573 results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- 574  
575 Peter L. Bartlett, Olivier Bousquet, and Shahar Mendelson. Local rademacher complexities. 2005.
- 576  
577 Peter L. Bartlett, Dylan J. Foster, and Matus J. Telgarsky. Spectrally-normalized margin bounds for  
578 neural networks. In *Advances in Neural Information Processing Systems*, pp. 6240–6249, 2017.
- 579  
580 Francesca Bartolucci, Ernesto De Vito, Lorenzo Rosasco, and Stefano Vigogna. Understanding  
581 neural networks with reproducing kernel banach spaces. *arXiv preprint arXiv:2109.09710*, 2021.
- 582  
583 Francesca Bartolucci, Ernesto De Vito, Lorenzo Rosasco, and Stefano Vigogna. Understanding  
584 neural networks with reproducing kernel banach spaces. *Applied and Computational Harmonic*  
585 *Analysis*, 62:194–236, January 2023.
- 586  
587 Brian Bell, Michaela Geyer, Juston Moore, David Glickenstein, and Amanda Fernandez. An exact  
588 kernel equivalence for finite classification models. In *TAG-ML*, 2023.
- 589  
590 John P. Boyd. The rate of convergence of hermite function series. *Mathematics of Computation*, 35  
591 (152):1309–1316, October 1980.
- 592  
593 Yuan Cao and Quanquan Gu. Generalization bounds of stochastic gradient descent for wide and deep  
neural networks. In *Advances in neural information processing systems*, volume 32, 2019.
- A. Caponnetto and E. De Vito. Optimal rates for the regularized least-squares algorithm. *Found*  
*Comput Math*, 7:331–368, 2007.
- C. Carmeli, E. De Vito, and A. Toigo. Reproducing kernel hilbert spaces and mercer theorem.  
Technical Report arXiv:math.FA/0504071, arXiv, 2005.

- 594 Youngmin Cho and Lawrence K. Saul. Kernel methods for deep learning. In Bengio Y., Schu-  
595 urmans D., Lafferty J. D., C. K. I. Williams, and A. Culotta (eds.), *Advances in Neural In-*  
596 *formation Processing Systems 22*, pp. 342–350. Curran Associates, Inc., 2009. URL [http:](http://papers.nips.cc/paper/3628-kernel-methods-for-deep-learning.pdf)  
597 [://papers.nips.cc/paper/3628-kernel-methods-for-deep-learning.pdf](http://papers.nips.cc/paper/3628-kernel-methods-for-deep-learning.pdf).
- 598 Corinna Cortes, Marius Kloft, and Mehryar Mohri. Learning kernels using local rademacher  
599 complexity. *Advances in neural information processing systems*, 26, 2013.
- 600 R. Courant and D. Hilbert. *Methods of Mathematical Physics*. John Wiley and sons, New York, 1937.
- 601 Amit Daniely. Sgd learns the conjugate kernel class of the network. In I. Guyon, U. V.  
602 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Ad-*  
603 *vances in Neural Information Processing Systems 30*, pp. 2422–2430. Curran Associates, Inc.,  
604 2017. URL [http://papers.nips.cc/paper/6836-sgd-learns-the-conjugate-](http://papers.nips.cc/paper/6836-sgd-learns-the-conjugate-kernel-class-of-the-network.pdf)  
605 [kernel-class-of-the-network.pdf](http://papers.nips.cc/paper/6836-sgd-learns-the-conjugate-kernel-class-of-the-network.pdf).
- 606 Amit Daniely, Roy Frostig, and Yoram Singer. Toward deeper understanding of neural networks:  
607 The power of initialization and a dual view on expressivity. In D. D. Lee, M. Sugiyama, U. V.  
608 Luxburg, I. Guyon, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*  
609 *29*, pp. 2253–2261. Curran Associates, Inc., 2016. URL [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/6427-toward-deeper-understanding-of-neural-networks-the-power-of-initialization-and-a-dual-view-on-expressivity.pdf)  
610 [6427-toward-deeper-understanding-of-neural-networks-the-power-](http://papers.nips.cc/paper/6427-toward-deeper-understanding-of-neural-networks-the-power-of-initialization-and-a-dual-view-on-expressivity.pdf)  
611 [of-initialization-and-a-dual-view-on-expressivity.pdf](http://papers.nips.cc/paper/6427-toward-deeper-understanding-of-neural-networks-the-power-of-initialization-and-a-dual-view-on-expressivity.pdf).
- 612 Ricky Der and Danial Lee. Large-margin classification in banach spaces. In *Proceedings of the JMLR*  
613 *Workshop and Conference 2: AISTATS2007*, pp. 91–98, 2007.
- 614 Felix Dräxler, Kambis Veschgini, Manfred Salmhofer, and Fred A. Hamprecht. Essentially no  
615 barriers in neural network energy landscape. In *Proceedings of the 35th International Conference*  
616 *on Machine Learning, ICML 2018*, 2018.
- 617 Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global  
618 minima of deep neural networks. In *International conference on machine learning*, pp. 1675–1685.  
619 PMLR, 2019a.
- 620 Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes  
621 over-parameterized neural networks. In *Conference on Learning Representations*, 2019b.
- 622 Adria Garriga-Alonso, Carl E. Rasmussen, and Laurence Aitchison. Deep convolutional networks as  
623 shallow gaussian processes. In *International Conference on Learning Representations*, pp. 1–16,  
624 May 2019.
- 625 Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In  
626 *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pp.  
627 315–323. JMLR Workshop and Conference Proceedings, 2011.
- 628 Noah Golowich, Alexander Rakhlin, and Ohad Shamir. Size-independent sample complexity of  
629 neural networks. In *COLT*, 2018.
- 630 Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. [http:](http://www.deeplearningbook.org)  
631 [://www.deeplearningbook.org](http://www.deeplearningbook.org).
- 632 I. S. Gradshteyn and I. M. Ryzhik. *Table of Integrals, Series, and Products*. Academic Press, London,  
633 2000.
- 634 Nick Harvey, Christopher Liaw, and Abbas Mehrabian. Nearly-tight VC-dimension bounds for  
635 piecewise linear neural networks. In *Proceedings of the 30th Conference on Learning Theory,*  
636 *COLT 2017*, 2017.
- 637 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
638 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
639 pp. 770–778, 2016.
- 640 Einar Hille. Contributions to the theory of Hermitian series. II. The representation problem. *Trans.*  
641 *Amer. Math. Soc.*, 47:80–94, 1940.

- 648 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,  
649 and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021.  
650
- 651 Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and  
652 generalization in neural networks. In *Advances in neural information processing systems*, pp. 8571–  
653 8580, 2018.
- 654 B. Laurent and P. Massart. Adaptive estimation of a quadratic functional by model selection. *The*  
655 *Annals of Statistics*, 28(5):1302 – 1338, 2000.  
656
- 657 Jaehoon Lee, Jascha Sohl-dickstein, Jeffrey Pennington, Roman Novak, Sam Schoenholz, and  
658 Yasaman Bahri. Deep neural networks as gaussian processes. In *In International Conference on*  
659 *Learning Representations*, 2018.  
660
- 661 Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-  
662 Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models  
663 under gradient descent. *Advances in neural information processing systems*, 32, 2019.
- 664 Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient  
665 descent on structured data. In *Advances in Neural Information Processing Systems 31: Annual*  
666 *Conference on Neural Information Processing Systems*, 2018.  
667
- 668 Rongrong Lin, Haizhang Zhang, and Jun Zhang. On reproducing kernel banach spaces: Generic  
669 definitions and unified framework of constructions. *Acta Mathematica Sinica, English Series*, 2022.  
670
- 671 Alexander G. de G. Matthews, Mark Rowland, Jiri Hron, Richard E. Turner, and Zoubin Ghahramani.  
672 Gaussian process behaviour in wide deep neural networks. *arXiv e-prints*, 2018.
- 673 James Mercer. Functions of positive and negative type, and their connection with the theory of  
674 integral equations. *Transactions of the Royal Society of London*, 209(A), 1909.  
675
- 676 Charles A. Micchelli and Massimiliano Pontil. On learning vector-valued functions. *Neural computa-*  
677 *tion*, 17(1):177–204, 2005.  
678
- 679 Philip M. Morse and Herman Feshbach. *Methods of Theoretical Physics*. McGraw-Hill, 1953.
- 680 Vaishnavh Nagarajan and J. Zico Kolter. Uniform convergence may be unable to explain gener-  
681 alization in deep learning. In H. Wallach, H. Larochelle, A. Beygelzimer, F. de S. Bengio, E. Fox,  
682 and R. Garnett (eds.), *Advances in Neural Information Processing Systems 32*, pp. 11615–  
683 11626. Curran Associates, Inc., 2019a. URL [http://papers.nips.cc/paper/9336-](http://papers.nips.cc/paper/9336-uniform-convergence-may-be-unable-to-explain-generalization-in-deep-learning.pdf)  
684 [uniform-convergence-may-be-unable-to-explain-generalization-in-](http://papers.nips.cc/paper/9336-uniform-convergence-may-be-unable-to-explain-generalization-in-deep-learning.pdf)  
685 [deep-learning.pdf](http://papers.nips.cc/paper/9336-uniform-convergence-may-be-unable-to-explain-generalization-in-deep-learning.pdf).
- 686 Vaishnavh Nagarajan and Zico Kolter. Deterministic PAC-Bayesian generalization bounds for deep  
687 networks via generalizing noise-resilience. In *International Conference on Learning Representa-*  
688 *tions (ICLR)*, 2019b.  
689
- 690 Radford M. Neal. *Priors for infinite networks*, pp. 29–53. Springer, 1996.  
691
- 692 Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. Norm-based capacity control in neural  
693 networks. In *Proceedings of Conference on Learning Theory*, pp. 1376–1401, 2015.
- 694 Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nathan Srebro. Exploring general-  
695 ization in deep learning. In *Proceedings of the 31st International Conference on Neural Informa-*  
696 *tion Processing Systems*, pp. 5949–5958, 2017.  
697
- 698 Behnam Neyshabur, Srinadh Bhojanapalli, and Nathan Srebro. A PAC-bayesian approach to  
699 spectrally-normalized margin bounds for neural networks. In *Proceedings of ICLR*, 2018.  
700
- 701 Behnam Neyshabur, Zhiyuan Li, Srinadh Bhojanapalli, Yann LeCun, and Nathan Srebro. The role of  
over-parametrization in generalization of neural networks. In *Proceedings of ICLR*, 2019.

- 702 Roman Novak, Lechao Xiao, Jaehoon Lee, Yasaman Bahri, Greg Yang, Jiri Hron, Dan Abolafia,  
703 Jeffrey Pennington, and Jascha Sohl-Dickstein. Bayesian deep convolutional networks with many  
704 channels are gaussian processes. In *International Conference on Learning Representations, ICLR*  
705 *2019*, 2019.
- 706 Frank W. Olver, Daniel W. Lozier, Ronald F. Boisvert, and Charles W. Clark. *NIST Handbook of*  
707 *Mathematical Functions*. Cambridge University Press, USA, 1st edition, 2010. ISBN 0521140633.
- 708 Rahul Parhi and Robert D. Nowak. Banach space representer theorems for neural networks and ridge  
709 splines. *J. Mach. Learn. Res.*, 22(43):1–40, 2021.
- 710 Ali Rahimi and Recht Benjamin. Weighted sums of random kitchen sinks: Replacing minimization  
711 with randomization in learning. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou (eds.),  
712 *Advances in Neural Information Processing Systems 21*, pp. 1313–1320. Curran Associates, Inc.,  
713 2009.
- 714 Marco Reiser and Hans Burkhardt. Learning equivariant functions with matrix valued kernels.  
715 *Journal of Machine Learning Research*, 8(15):385–408, 2007.
- 716 Koen Sanders. Neural networks as functions parameterized by measures: Representer theorems and  
717 approximation benefits. Master’s thesis, Eindhoven University of Technology, 2020.
- 718 Laurent Schwartz. Sous-espaces hilbertiens d’espaces vectoriels topologiques et noyaux associés  
719 (noyaux reproduisants). *Journal d’analyse mathématique*, 13:115–256, 1964.
- 720 John Shawe-Taylor and Nello Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University  
721 Press, 2004.
- 722 Alistair Shilton, Sunil Gupta, Santu Rana, and Svetha Venkatesh. Gradient descent in neural networks  
723 as sequential learning in reproducing kernel banach space. In Andreas Krause, Emma Brunskill,  
724 Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of*  
725 *the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine*  
726 *Learning Research*, pp. 31435–31488. PMLR, 23–29 Jul 2023.
- 727 Guohui Song, Haizhang Zhang, and Fred J. Hickernell. Reproducing kernel banach spaces with the  
728  $\ell^1$  norm. *Applied and Computational Harmonic Analysis*, 34(1):96–116, Jan 2013.
- 729 Len Spek, Tjeerd Jan Heeringa, and Christoph Brune. Duality for neural networks through reproduc-  
730 ing kernel banach spaces. *arXiv preprint arXiv:2211.05020*, 2022.
- 731 Bharath K. Sriperumbudur, Kenji Fukumizu, and Gert R. Lanckriet. Learning in hilbert vs. banach  
732 spaces: A measure embedding viewpoint. In *Advances in Neural Information Processing Systems*,  
733 pp. 1773–1781, 2011.
- 734 Ingo Steinwart and Andreas Christman. *Support Vector Machines*. Springer, 2008.
- 735 Michael Unser. A representer theorem for deep neural networks. *J. Mach. Learn. Res.*, 20(110):1–30,  
736 2019.
- 737 Michael Unser. A unifying representer theorem for inverse problems and machine learning. *Founda-*  
738 *tions of Computational Mathematics*, 21(4):941–960, 2021.
- 739 E. Weinan, Chao Ma, and Lei Wu. A priori estimates of the population risk for two-layer neural  
740 networks. *Communications in Mathematical Sciences*, 17(5):1407–1425, 2019.
- 741 Yuesheng Xu and Qi Ye. Generalized mercer kernels and reproducing kernel banach spaces. *arXiv*  
742 *preprint arXiv:1412.8663*, 2014.
- 743 Haizhang Zhang and Jun Zhang. Regularized learning in banach spaces as an optimization problem:  
744 representer theorems. *Journal of Global Optimization*, 54(2):235–250, Oct 2012.
- 745 Haizhang Zhang, Yuesheng Xu, and Jun Zhang. Reproducing kernel banach spaces for machine  
746 learning. *Journal of Machine Learning Research*, 10:2741–2775, 2009.

756 Wenda Zhou, Victor Veitch, Morgane Austern, Ryan P. Adams, and Peter Orbanz. Nonvacuous gen-  
757 eralization bounds at the imagenet scale: a PAC-Bayesian compression approach. In *International*  
758 *Conference on Learning Representations (ICLR)*, 2019.

759 Difan Zou and Quanquan Gu. An improved analysis of training over-parameterized deep neural  
760 networks. In *Advances in neural information processing systems*, volume 32, 2019.

761 Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Gradient descent optimizes over-  
762 parameterized deep relu networks. *Machine learning*, 109(3):467–492, 2020.

763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

## A PROPERTIES OF HERMITE POLYNOMIALS

The (probabilist's) Hermite polynomials are given by (Abramowitz et al., 1972; Morse & Feshbach, 1953; Olver et al., 2010; Courant & Hilbert, 1937):

$$He_k(\zeta) = (-1)^k e^{\frac{\zeta^2}{2}} \frac{d^k}{d\zeta^k} e^{-\frac{\zeta^2}{2}} \quad \forall k \in \mathbb{N}$$

or, explicitly:

$$He_k(\zeta) = k! \sum_{0 \leq 2p \leq k} \frac{(-1)^p}{2^p p! (k-2p)!} \zeta^{k-2p} \quad \forall k \in \mathbb{N} \quad (24)$$

and form an orthogonal basis of  $L^2(\mathbb{R}, e^{-x^2})$ . Thus for any  $f \in L^2(\mathbb{R}, e^{-x^2})$  there exist Hermite coefficients  $a_0, a_1, \dots \in \mathbb{R}$  (i.e. the Hermite transform of  $f$ ) so that:

$$f(\zeta) = \sum_{k \in \mathbb{N}} a_k He_k(\zeta) \quad \forall \zeta \in \mathbb{R}$$

where:

$$a_k = \frac{1}{k! \sqrt{2\pi}} \int_{-\infty}^{\infty} f(\zeta) e^{-\frac{\zeta^2}{2}} He_k(\zeta) d\zeta$$

This series representation converges everywhere on the real line. Moreover (Hille, 1940; Boyd, 1980) this series converges on a strip  $\mathbb{S}_\rho = \{z \in \mathbb{C} : -\rho < \text{Im}(z) < \rho\}$  of width  $\rho$  about the real axis in the complex plane, where:<sup>7</sup>

$$\rho = -\limsup_{k \rightarrow \infty} \frac{1}{\sqrt{2k+1}} \log \left( \left| \frac{a_k}{\sqrt{k! \sqrt{\pi}}} \right| \right) \quad (25)$$

The Hermite numbers derive from the Hermite polynomials:<sup>8</sup>

$$He_k \triangleq He_k(0) = \begin{cases} 0 & \text{if } k \text{ odd} \\ \frac{k!}{(\frac{k}{2})!} \left(-\frac{1}{2}\right)^{\frac{k}{2}} & \text{if } k \text{ even} \end{cases}$$

It is well known that (see e.g. (Morse & Feshbach, 1953)):

$$He_k(\zeta + \xi) = \sum_{0 \leq l \leq k} \binom{k}{l} He_{k-l}(\zeta) \xi^l$$

and so:

$$He_k(\zeta) = \sum_{0 \leq l \leq k} \binom{k}{l} He_{k-l} \zeta^l$$

It follows that, taking care not to change or order of summation (remember this is an alternating series):

$$f(\zeta) = \sum_{k=0}^{\infty} a_k \sum_{l=0}^k \binom{k}{l} He_{k-l} \zeta^l$$

Next we derive a helpful property involving *rectified* Hermite expansions. Let  $f \in L^2(\mathbb{R}, e^{-x^2})$ ,  $f(0) = 0$ , then:

$$f(x) = \sum_{k=1}^{\infty} a_k He_k(x) = \sum_{k=1}^{\infty} a_k \sum_{l=1}^k \binom{k}{l} He_{k-l} x^l$$

Denoting the imaginary element  $i$ :

$$\begin{aligned} f(ix) &= \sum_{k=1}^{\infty} a_k \sum_{l=1}^k \binom{k}{l} He_{k-l}(ix)^l \\ &= \sum_{k=1}^{\infty} a_k \sum_{l=0}^{k-1} \binom{k}{k-l} He_l(ix)^{k-l} \\ &= \sum_{k=1}^{\infty} i^k a_k \sum_{l=0}^{k-1} \binom{k}{k-l} (i^l He_l) x^{k-l} \end{aligned}$$

<sup>7</sup>Note that (Hille, 1940; Boyd, 1980) use the normalized physicist's Hermite polynomials. The additional scale factor here arises in the translation to the un-normalized probabilist's Hermite polynomials used here.

<sup>8</sup>Typically the Hermite numbers are defined from the physicist's Hermite polynomials, but as we use the Probabilist's form we find these more convenient.

864 Recall that  $\text{He}_l = 0$  for  $l = 1, 3, 5, \dots$ , and  $\text{sgn}(\text{He}_{2p}) = (-1)^p$ . Therefore

$$\begin{aligned} 866 f(\mathbf{i}x) &= \sum_{k=1}^{\infty} \mathbf{i}^k a_k \sum_{l=0}^{k-1} \binom{k}{k-l} |\text{He}_l| x^{k-l} \\ 867 &= \sum_{k=1}^{\infty} \mathbf{i}^k a_k \sum_{l=1}^k \binom{k}{l} |\text{He}_{k-l}| x^l \end{aligned}$$

871 and so:

$$\begin{aligned} 873 \text{Im}(f(\mathbf{i}x)) &= \sum_{k=1,3,5,\dots} (-1)^{\lfloor \frac{k}{2} \rfloor} a_k \sum_{l=1}^k \binom{k}{l} |\text{He}_{k-l}| x^l \\ 874 \\ 875 \text{Re}(f(\mathbf{i}x)) &= \sum_{k=2,4,6,\dots} (-1)^{\lfloor \frac{k}{2} \rfloor} a_k \sum_{l=1}^k \binom{k}{l} |\text{He}_{k-l}| x^l \end{aligned}$$

$$877 \text{Im}(f(\mathbf{i}x)) + \text{Re}(f(\mathbf{i}x)) = \sum_{k=1}^{\infty} (-1)^{\lfloor \frac{k}{2} \rfloor} a_k \sum_{l=1}^k \binom{k}{l} |\text{He}_{k-l}| x^l \quad (26)$$

882 Finally we make some observations regarding derivatives that will be required later. Let  $f \in$   
883  $L^2(\mathbb{R}, e^{-x^2})$ ,  $f(0) = 0$ . Then:

$$885 f(x) = \sum_{k=1}^{\infty} a_k \sum_{l=1}^k \binom{k}{l} \text{He}_{k-l} x^l$$

888 and subsequently:

$$\begin{aligned} 890 f^{(1)}(x) &= \sum_{k=1}^{\infty} a_k \sum_{l=1}^k l \binom{k}{l} \text{He}_{k-l} x^{l-1} \\ 891 &= \sum_{k=1}^{\infty} a_k \sum_{l=1}^k l \frac{k!}{l!(k-l)!} \text{He}_{k-l} x^{l-1} \\ 892 &= \sum_{k=1}^{\infty} a_k \sum_{l=1}^k k \frac{(k-1)!}{(l-1)!(k-l)!} \text{He}_{k-l} x^{l-1} \\ 893 &= \sum_{k=1}^{\infty} a_k \sum_{l=1}^k k \binom{k-1}{l-1} \text{He}_{k-l} x^{l-1} \\ 894 &= \sum_{k=0}^{\infty} (k+1) a_{k+1} \sum_{l=0}^{k-1} \binom{k}{l} \text{He}_{k-l} x^l \end{aligned}$$

901 and:

$$\begin{aligned} 902 f^{(2)}(x) &= \sum_{k=0}^{\infty} (k+1) a_{k+1} \sum_{l=1}^{k-1} (l-1) \binom{k}{l} \text{He}_{k-l} x^{l-1} \\ 903 &= \sum_{k=0}^{\infty} (k+1) a_{k+1} \sum_{l=1}^{k-1} (l-1) \frac{k!}{l!(k-l)!} \text{He}_{k-l} x^{l-1} \\ 904 &= \sum_{k=0}^{\infty} (k+1) a_{k+1} \sum_{l=1}^{k-1} k \frac{(k-1)!}{(l-1)!(k-l)!} \text{He}_{k-l} x^{l-1} \\ 905 &= \sum_{k=1}^{\infty} (k+1) a_{k+1} \sum_{l=1}^{k-1} k \frac{(k-1)!}{(l-1)!(k-l)!} \text{He}_{k-l} x^{l-1} \\ 906 &= \sum_{k=0}^{\infty} (k+1)(k+2) a_{k+2} \sum_{l=0}^{k-2} \binom{k}{l} \text{He}_{k-l} x^l \end{aligned}$$

914 and so on to:

$$915 f^{(n)}(x) = \sum_{k=0}^{\infty} \frac{(k+n)!}{k!} a_{k+n} \sum_{l=0}^{k-n} \binom{k}{l} \text{He}_{k-l} x^l$$

## A.1 ACTIVATION FUNCTIONS

Following the previous method we introduce our notation for the activation functions. Recall  $\tau^{[\tilde{j},j]} \in L^2(\mathbb{R}, e^{-x^2})$  by assumption. Subsequently  $\bar{\tau}^{[\tilde{j},j]} \in L^2(\mathbb{R}, e^{-x^2})$ , where:

$$\begin{aligned} \bar{\tau}^{[\tilde{j},j]}(\zeta; \xi) &= \tau^{[\tilde{j},j]}(\xi + \zeta) - \tau^{[\tilde{j},j]}(\xi) = \sum_{k \in \mathbb{N}} a_{(\xi)k}^{[\tilde{j},j]} \text{He}_k(\zeta) \quad \forall \zeta \in \mathbb{R}_+ \\ &= \sum_{k=0}^{\infty} a_{(\xi)k}^{[\tilde{j},j]} \sum_{l=0}^k \binom{k}{l} \text{He}_{k-l} \zeta^l \\ &= \sum_{k=1}^{\infty} a_{(\xi)k}^{[\tilde{j},j]} \sum_{l=1}^k \binom{k}{l} \text{He}_{k-l} \zeta^l \end{aligned} \quad (27)$$

(in the final step we use that  $\bar{\tau}^{[\tilde{j},j]}(0; \xi) = 0$ ) with coefficients:

$$a_{(\xi)k}^{[\tilde{j},j]} = \frac{1}{k! \sqrt{2\pi}} \int_{-\infty}^{\infty} \bar{\tau}^{[\tilde{j},j]}(\zeta; \xi) e^{-\frac{\zeta^2}{2}} \text{He}_k(\zeta) d\zeta$$

which converges on a strip  $\mathbb{S}_{\rho_{(\xi)}^{[\tilde{j},j]}} = \{z \in \mathbb{C} : -\rho_{(\xi)}^{[\tilde{j},j]} < \text{Im}(z) < \rho_{(\xi)}^{[\tilde{j},j]}\}$  of width  $\rho_{(\xi)}^{[\tilde{j},j]}$  about the real axis in the complex plane, where:

$$\rho_{(\xi)}^{[\tilde{j},j]} = -\limsup_{k \rightarrow \infty} \frac{1}{\sqrt{2k+1}} \log \left( \left| \frac{a_{(\xi)k}^{[\tilde{j},j]}}{\sqrt{k! \sqrt{\pi}}} \right| \right)$$

## A.2 RECTIFIED ACTIVATION FUNCTIONS

Recall that the rectified activation functions are defined as:

$$\hat{\tau}_{\eta}^{[\tilde{j},j]}(\zeta; \xi, \xi') = \sum_{k=1}^{\infty} \frac{a_{(\xi)k}^{[\tilde{j},j]} a_{(\xi')k}^{[\tilde{j},j]}}{\eta^{2k}} \sum_{l=1}^k \binom{k}{l}^2 \text{He}_{k-l}^2 \zeta^l$$

where  $\eta \in (0, 1)$  is fixed. To understand the convergence of this function, observe that:

$$\hat{\tau}_{\eta}^{[\tilde{j},j]}(\zeta; \xi, \xi') \leq \max_{\xi'' \in \{\xi, \xi'\}} \sum_{k=1}^{\infty} \sum_{l=1}^k \left| \frac{a_{(\xi'')k}^{[\tilde{j},j]}}{\eta^k} \binom{k}{l} \text{He}_{k-l} \sqrt{\zeta}^l \right|^2$$

which is the 2-norm of a sequence. Hence:

$$\hat{\tau}_{\eta}^{[\tilde{j},j]}(\zeta; \xi, \xi') \leq \max_{\xi'' \in \{\xi, \xi'\}} \left( \sum_{k=1}^{\infty} \sum_{l=1}^k \left| \frac{a_{(\xi'')k}^{[\tilde{j},j]}}{\eta^k} \binom{k}{l} |\text{He}_{k-l}| \sqrt{\zeta}^l \right|^2 \right)^2$$

Thus it suffices to study the convergence of:

$$\gamma^{[\tilde{j},j]}(\lambda; \xi) = \sum_{k=1}^{\infty} \frac{|a_{(\xi)k}^{[\tilde{j},j]}|}{\eta^k} \sum_{l=1}^k \binom{k}{l} |\text{He}_{k-l}| \lambda^l$$

which in turn bounds:

$$\hat{\tau}_{\eta}^{[\tilde{j},j]}(\zeta; \xi, \xi') \leq \max \left\{ \gamma^{[\tilde{j},j]}(\sqrt{\zeta}; \xi)^2, \gamma^{[\tilde{j},j]}(\sqrt{\zeta}; \xi')^2 \right\}$$

Using (26):

$$\gamma^{[\tilde{j},j]}(\lambda; \xi) = \text{Re}(\bar{\gamma}^{[\tilde{j},j]}(i\lambda; \xi)) + \text{Im}(\bar{\gamma}^{[\tilde{j},j]}(i\lambda; \xi))$$

where:

$$\gamma^{[\tilde{j},j]}(\lambda; \xi) = \sum_{k=1}^{\infty} \frac{|a_{(\xi)k}^{[\tilde{j},j]}|}{\eta^k} \sum_{l=1}^k \binom{k}{l} \text{He}_{k-l} \lambda^l$$

which, using (25), is convergent for:

$$|\lambda| < L_{(\xi)}^{[\tilde{j},j]} < -\limsup_{k \rightarrow \infty} \frac{1}{\sqrt{2k+1}} \left( \ln \left| \frac{a_{(\xi)k}^{[\tilde{j},j]}}{\eta^k} \right| \right)$$

whereon:

$$|\gamma^{[\tilde{j},j]}(\lambda; \xi)| \leq M_{(\xi)}^{[\tilde{j},j]} = \gamma^{[\tilde{j},j]}(L_{(\xi)}^{[\tilde{j},j]}; \xi)$$

and we conclude that  $\hat{\tau}_\eta^{[\bar{j},j]}(\zeta; \xi, \xi')$  is convergent for:

$$|\zeta| \leq \rho_{(\xi)\eta}^{[\bar{j},j]2} = \min \left\{ -\limsup_{k \rightarrow \infty} \frac{1}{\sqrt{2k+1}} \ln \left| \frac{a_{(\xi)k}^{[\bar{j},j]}}{\eta^k} \right|, -\limsup_{k \rightarrow \infty} \frac{1}{\sqrt{2k+1}} \ln \left| \frac{a_{(\xi)k}^{[\bar{j},j]}}{\eta^k} \right| \right\}^2$$

whereon:

$$\hat{\tau}_\eta^{[\bar{j},j]}(\zeta; \xi, \xi') \leq T_{(\xi)\eta}^{[\bar{j},j]2} = \hat{\tau}_\eta^{[\bar{j},j]}(\rho_{(\xi)\eta}^{[\bar{j},j]2}; \xi, \xi')$$

The envelope is convergent for:

$$|\zeta| < \rho_{(\omega)\eta}^{[\bar{j},j]2} = \inf_{|\xi|, |\xi'| \leq \omega} \rho_{(\xi, \xi')\eta}^{[\bar{j},j]2}$$

whereon:

$$T_{(\omega)\eta}^{[\bar{j},j]2} = \hat{\tau}_\eta^{[\bar{j},j]}(\rho_{(\omega)\eta}^{[\bar{j},j]2}; \omega)$$

Finally:

$$\begin{aligned} \hat{s}_\eta(\zeta) &= \sum_{k=1}^{\infty} \eta^{2k} \sum_{l=1}^k \zeta^l = \frac{\zeta}{1-\zeta} \sum_{k=1}^{\infty} \eta^{2k} (1 - \zeta^k) = \frac{\zeta}{1-\zeta} \left( \sum_{k=1}^{\infty} \eta^{2k} - \sum_{k=1}^{\infty} (\zeta \eta^2)^k \right) \\ &= \frac{\zeta}{1-\zeta} \left( \frac{\eta^2}{1-\eta^2} - \frac{\zeta \eta^2}{1-\zeta \eta^2} \right) \end{aligned}$$

is convergent  $\forall |\zeta| < R_\eta^2 < 1$ , with max value of  $S_\eta^2 = \frac{R_\eta^2}{1-R_\eta^2} \left( \frac{\eta^2}{1-\eta^2} - \frac{\eta^2 R_\eta^2}{1-\eta^2 R_\eta^2} \right)$  thereon.

### A.3 RELU ACTIVATION FUNCTION

In this section we derive the Hermite-polynomial expansion of the centered ReLU activation function:

$$\begin{aligned} \bar{\tau}^{[\text{ReLU}]}(\zeta; \xi) &= \tau^{[\text{ReLU}]}(\xi + \zeta) - \tau^{[\text{ReLU}]}(\xi) \\ &= \begin{cases} \zeta + \xi & \text{if } \zeta > -\xi \\ 0 & \text{otherwise} \end{cases} - [\xi]_+ \\ &= \sum_{k=0}^{\infty} a_{(\xi)k}^{[\bar{j},j]} He_k(\zeta) \end{aligned}$$

We find it convenient to work in terms of the physicists Hermite polynomials  $H_k$  to suit (Gradshteyn & Ryzhik, 2000). Using this:

$$\begin{aligned} a_{(\xi)k}^{[\text{ReLU}]} &= \frac{1}{\sqrt{2\pi k!}} \int_{-\xi}^{\infty} (\zeta + \xi) e^{-\frac{\zeta^2}{2}} He_k(\zeta) d\zeta - \frac{1}{\sqrt{2\pi k!}} \int_{-\infty}^{\infty} [\xi]_+ e^{-\frac{\zeta^2}{2}} He_k(\zeta) d\zeta \\ &= \frac{1}{\sqrt{2\pi k!}} \int_{-\xi}^{\infty} (\zeta + \xi) e^{-\frac{\zeta^2}{2}} \frac{1}{\sqrt{2^k}} H_k \left( \frac{\zeta}{\sqrt{2}} \right) d\zeta - \frac{1}{\sqrt{2\pi k!}} \int_{-\infty}^{\infty} [\xi]_+ e^{-\frac{\zeta^2}{2}} \frac{1}{\sqrt{2^k}} H_k \left( \frac{\zeta}{\sqrt{2}} \right) d\zeta \\ &= \sqrt{\frac{2}{\pi}} \frac{1}{k!} \int_{-\frac{\xi}{\sqrt{2}}}^{\infty} \left( \frac{\zeta}{\sqrt{2}} + \frac{\xi}{\sqrt{2}} \right) e^{-\left(\frac{\zeta}{\sqrt{2}}\right)^2} \frac{1}{\sqrt{2^k}} H_k \left( \frac{\zeta}{\sqrt{2}} \right) d\zeta - \sqrt{\frac{2}{\pi}} \frac{1}{k!} \int_{-\infty}^{\infty} \left[ \frac{\xi}{\sqrt{2}} \right]_+ e^{-\left(\frac{\zeta}{\sqrt{2}}\right)^2} \frac{1}{\sqrt{2^k}} H_k \left( \frac{\zeta}{\sqrt{2}} \right) d\zeta \\ &= \sqrt{\frac{2}{\pi}} \frac{1}{\sqrt{2^k k!}} \int_{-\frac{\xi}{\sqrt{2}}}^{\infty} (\zeta + \frac{\xi}{\sqrt{2}}) e^{-\zeta^2} H_k(\zeta) d\zeta - \sqrt{\frac{2}{\pi}} \frac{1}{\sqrt{2^k k!}} \int_{-\infty}^{\infty} \left[ \frac{\xi}{\sqrt{2}} \right]_+ e^{-\zeta^2} H_k(\zeta) d\zeta \\ &= \sqrt{\frac{2}{\pi}} \frac{1}{\sqrt{2^k k!}} \int_{-\frac{\xi}{\sqrt{2}}}^{\infty} e^{-\zeta^2} \zeta H_k(\zeta) d\zeta + \frac{1}{\sqrt{\pi}} \frac{1}{\sqrt{2^k k!}} \xi \int_{-\frac{\xi}{\sqrt{2}}}^{\infty} e^{-\zeta^2} H_k(\zeta) d\zeta - \frac{1}{\sqrt{\pi}} \frac{1}{\sqrt{2^k k!}} [\xi]_+ \int_{-\infty}^{\infty} e^{-\zeta^2} H_k(\zeta) d\zeta \end{aligned}$$

Using the recursion and derivative properties, for  $k > 1$ :

$$\begin{aligned} \zeta H_k(\zeta) &= \frac{1}{2} H_{k+1}(\zeta) + \frac{1}{2} H'_k(\zeta) \\ &= \frac{1}{2} H_{k+1}(\zeta) + k H_{k-1}(\zeta) \end{aligned}$$

and hence, using (Gradshteyn & Ryzhik, 2000, (7.373)):

$$\begin{aligned} a_{(\xi)k}^{[\text{ReLU}]} &= \frac{k+1}{\sqrt{\pi}} \frac{1}{\sqrt{2^{k+1}} (k+1)!} \int_{-\frac{\xi}{\sqrt{2}}}^{\infty} e^{-\zeta^2} H_{k+1}(\zeta) d\zeta + \frac{1}{\sqrt{\pi}} \frac{1}{\sqrt{2^{k-1}} (k-1)!} \int_{-\frac{\xi}{\sqrt{2}}}^{\infty} e^{-\zeta^2} H_{k-1}(\zeta) d\zeta + \dots \\ &\quad \dots + \frac{1}{\sqrt{\pi}} \frac{1}{\sqrt{2^k k!}} \xi \int_{-\frac{\xi}{\sqrt{2}}}^{\infty} e^{-\zeta^2} H_k(\zeta) d\zeta - \frac{1}{\sqrt{\pi}} \frac{1}{\sqrt{2^k k!}} [\xi]_+ \int_{-\infty}^{\infty} e^{-\zeta^2} H_k(\zeta) d\zeta \\ &= \frac{k+1}{\sqrt{\pi}} \frac{1}{\sqrt{2^{k+1}} (k+1)!} \int_{-\frac{\xi}{\sqrt{2}}}^{\infty} e^{-\zeta^2} H_{k+1}(\zeta) d\zeta + \frac{1}{\sqrt{\pi}} \frac{1}{\sqrt{2^{k-1}} (k-1)!} \int_{-\frac{\xi}{\sqrt{2}}}^{\infty} e^{-\zeta^2} H_{k-1}(\zeta) d\zeta + \dots \\ &\quad \dots + \frac{1}{\sqrt{\pi}} \frac{1}{\sqrt{2^k k!}} \xi \int_{-\frac{\xi}{\sqrt{2}}}^{\infty} e^{-\zeta^2} H_k(\zeta) d\zeta - \delta_{k,2p} \frac{1}{\sqrt{\pi}} \frac{1}{\sqrt{2^{2p}} (2p)!} [\xi]_+ \int_{-\infty}^{\infty} e^{-\zeta^2} H_{2p}(\zeta) d\zeta \\ &= \frac{k+1}{\sqrt{\pi}} \frac{1}{\sqrt{2^{k+1}} (k+1)!} \int_{-\frac{\xi}{\sqrt{2}}}^{\infty} e^{-\zeta^2} H_{k+1}(\zeta) d\zeta + \frac{1}{\sqrt{\pi}} \frac{1}{\sqrt{2^{k-1}} (k-1)!} \int_{-\frac{\xi}{\sqrt{2}}}^{\infty} e^{-\zeta^2} H_{k-1}(\zeta) d\zeta + \dots \\ &\quad \dots + \frac{1}{\sqrt{\pi}} \frac{1}{\sqrt{2^k k!}} \xi \int_{-\frac{\xi}{\sqrt{2}}}^{\infty} e^{-\zeta^2} H_k(\zeta) d\zeta \end{aligned}$$

1026

Using (Gradshteyn &amp; Ryzhik, 2000, (7.373)) we have:

1027

1028

1029

1030

1031

1032

1033

1034

1035

$$\begin{aligned}
a_{(\xi)k}^{[\text{ReLU}]} &= \frac{1}{\sqrt{\pi}} \frac{1}{\sqrt{2^{k+1}}(k+1)!} (k+1) \left( e^{-\frac{\xi^2}{2}} H_k \left( -\frac{\xi}{\sqrt{2}} \right) - e^{-\frac{\infty^2}{2}} H_k \left( \frac{\infty}{\sqrt{2}} \right) \right) + \dots \\
&\quad \dots \frac{1}{\sqrt{\pi}} \frac{1}{\sqrt{2^{k-1}}(k-1)!} \left( e^{-\frac{\xi^2}{2}} H_{k-2} \left( -\frac{\xi}{\sqrt{2}} \right) - e^{-\frac{\infty^2}{2}} H_{k-2} \left( \frac{\infty}{\sqrt{2}} \right) \right) + \dots \\
&\quad \dots \frac{1}{\sqrt{\pi}} \frac{1}{\sqrt{2^k k!}} \xi \left( e^{-\frac{\xi^2}{2}} H_{k-1} \left( -\frac{\xi}{\sqrt{2}} \right) - e^{-\frac{\infty^2}{2}} H_{k-1} \left( \frac{\infty}{\sqrt{2}} \right) \right) \\
&= \frac{1}{\sqrt{2\pi}} e^{-\frac{\xi^2}{2}} \left( \frac{k+1}{\sqrt{2^k}(k+1)!} H_k \left( -\frac{\xi}{\sqrt{2}} \right) + \frac{1}{\sqrt{2^{k-2}}(k-1)!} H_{k-2} \left( -\frac{\xi}{\sqrt{2}} \right) + \frac{1}{\sqrt{2^{k-1}}k!} \xi H_{k-1} \left( -\frac{\xi}{\sqrt{2}} \right) \right)
\end{aligned}$$

If  $k = 2p$  and  $p > 0$  then, noting that  $H_k(0) = \sqrt{2}^k \text{He}_k$ :

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

If  $k = 2p + 1$  and  $p > 0$  then:

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

For the cases  $k = 0, 1$  we use the result:

$$\int_a^b x^m e^{-x^2} dx = \frac{1}{2} \Gamma \left( \frac{m+1}{2}, a^2 \right) - \frac{1}{2} \Gamma \left( \frac{m+1}{2}, b^2 \right)$$

and so:

$$\int_a^\infty x^m e^{-x^2} dx = \frac{1}{2} \Gamma \left( \frac{m+1}{2}, a^2 \right)$$

In the case  $k = 0$ :

$$\begin{aligned}
a_{(\xi)0}^{[\text{ReLU}]} &= \sqrt{\frac{2}{\pi}} \int_{-\frac{\xi}{\sqrt{2}}}^\infty \zeta e^{-\zeta^2} d\zeta + \frac{1}{\sqrt{\pi}} \xi \int_{-\frac{\xi}{\sqrt{2}}}^\infty e^{-\zeta^2} d\zeta - [\xi]_+ \\
&= \frac{1}{\sqrt{2\pi}} \Gamma \left( 1, \frac{\xi^2}{2} \right) + \frac{1}{2\sqrt{\pi}} \xi \Gamma \left( \frac{1}{2}, \frac{\xi^2}{2} \right) - [\xi]_+ \\
&= \frac{1}{\sqrt{2\pi}} \text{ if } \xi = 0
\end{aligned}$$

and in the case  $k = 1$ :

$$\begin{aligned}
a_{(\xi)1}^{[\text{ReLU}]} &= \frac{2}{\sqrt{\pi}} \int_{-\frac{\xi}{\sqrt{2}}}^\infty \zeta^2 e^{-\zeta^2} d\zeta + \sqrt{\frac{2}{\pi}} \xi \int_{-\frac{\xi}{\sqrt{2}}}^\infty \zeta e^{-\zeta^2} d\zeta \\
&= \frac{1}{\sqrt{\pi}} \Gamma \left( \frac{3}{2}, \frac{\xi^2}{2} \right) + \frac{1}{\sqrt{2\pi}} \xi \Gamma \left( 1, \frac{\xi^2}{2} \right) \\
&= \frac{1}{2} \text{ if } \xi = 0
\end{aligned}$$

Next we derive the magnitude functions for the ReLU. Using integration by parts, we see that:

1073

1074

1075

1076

1077

1078

1079

$$\begin{aligned}
\frac{1}{\sqrt{2\pi}} \int_c^x \frac{1}{\zeta^2} \left( e^{\frac{1}{2}\zeta^2} - 1 \right) d\zeta &= \frac{1}{\sqrt{2\pi}} \frac{1}{\sqrt{2}} \int_c^x \frac{2}{\zeta^2} \left( e^{\frac{1}{2}\zeta^2} - 1 \right) d\frac{\zeta}{\sqrt{2}} \\
&= \frac{1}{\sqrt{2\pi}} \frac{1}{\sqrt{2}} \int_{\frac{c}{\sqrt{2}}}^{\frac{x}{\sqrt{2}}} \frac{1}{\zeta^2} \left( e^{\zeta^2} - 1 \right) d\zeta \\
&= -\frac{1}{\sqrt{2\pi}} \frac{1}{x} \left( e^{\frac{1}{2}x^2} - 1 \right) + \frac{1}{\sqrt{2\pi}} \frac{1}{c} \left( e^{\frac{1}{2}c^2} - 1 \right) + \frac{1}{\sqrt{\pi}} \int_{\frac{c}{\sqrt{2}}}^{\frac{x}{\sqrt{2}}} e^{\zeta^2} d\zeta \\
&= -\frac{1}{\sqrt{2\pi}} \frac{1}{x} \left( e^{\frac{1}{2}x^2} - 1 \right) + \frac{1}{\sqrt{2\pi}} \frac{1}{c} \left( e^{\frac{1}{2}c^2} - 1 \right) + \frac{1}{2} \frac{2}{\sqrt{\pi}} \int_{\frac{c}{\sqrt{2}}}^{\frac{x}{\sqrt{2}}} e^{\zeta^2} d\zeta \\
&= -\frac{1}{\sqrt{2\pi}} \frac{1}{x} \left( e^{\frac{1}{2}x^2} - 1 \right) + \frac{1}{2} \text{erfi} \left( \frac{x}{\sqrt{2}} \right) - \frac{1}{2} \left( \text{erfi} \left( \frac{c}{\sqrt{2}} \right) - \frac{1}{\sqrt{2\pi}} \frac{2}{c} \left( e^{\frac{1}{2}c^2} - 1 \right) \right)
\end{aligned}$$

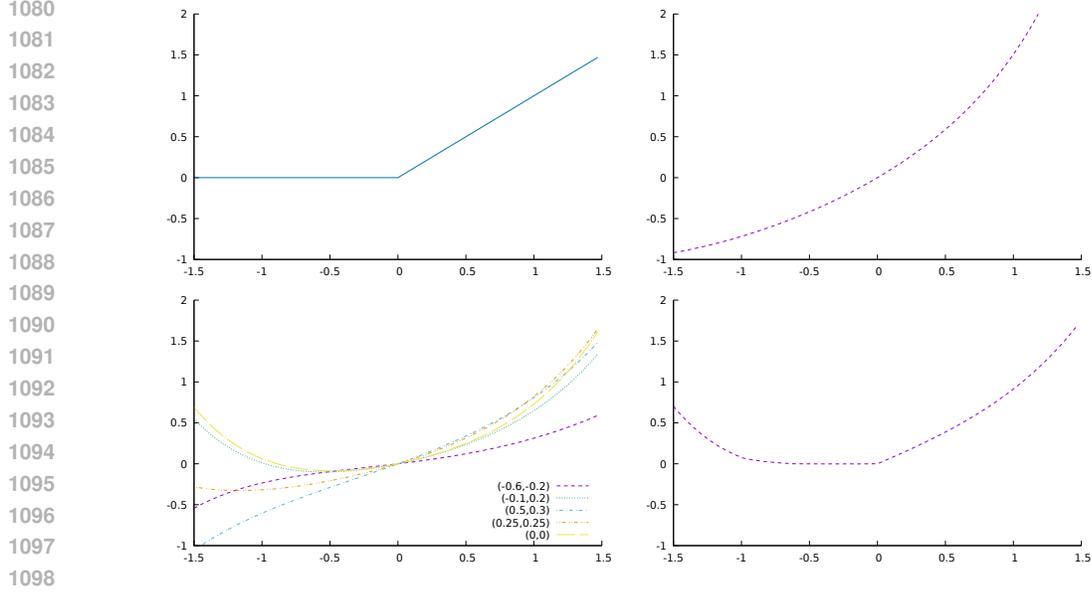


Figure 3: The ReLU magnitude activation and associated functions - ReLU activation  $\tau^{[\text{ReLU}]}$  (top left), magnitude  $\bar{s}^{[\text{ReLU}]}$  (top right), rectified activation  $\hat{\tau}_\eta^{[\text{ReLU}]}$  for various  $\xi, \xi'$  pairs (bottom left) and its envelope (bottom right) for  $\xi, \xi' \in [-2, 2]$ . Note that  $\eta = 0.75$  for all plots.

So:

$$\begin{aligned}
\sum_{k=1}^{\infty} |a_k^{[\text{ReLU}]}| x^k &= \frac{1}{2}x + \frac{1}{\sqrt{2\pi}} \sum_{p=1}^{\infty} \frac{x^{2p}}{(2p-1)2^p p!} \\
&= \frac{1}{2}x + \frac{1}{\sqrt{2\pi}} x \sum_{p=1}^{\infty} \frac{x^{2p-1}}{(2p-1)2^p p!} \\
&= \frac{1}{2}x + \frac{1}{\sqrt{2\pi}} x \int_c^x \left( \frac{\partial}{\partial \zeta} \sum_{p=1}^{\infty} \frac{\zeta^{2p-1}}{(2p-1)2^p p!} \right) d\zeta \\
&= \frac{1}{2}x + \frac{1}{\sqrt{2\pi}} x \int_c^x \left( \sum_{p=1}^{\infty} \frac{\zeta^{2p-2}}{2^p p!} \right) d\zeta \\
&= \frac{1}{2}x + \frac{1}{2\sqrt{2\pi}} x \int_c^x \left( \sum_{p=1}^{\infty} \frac{1}{p!} \left(\frac{1}{2}\zeta^2\right)^{p-1} \right) d\zeta \\
&= \frac{1}{2}x + \frac{1}{\sqrt{2\pi}} x \int_c^x \frac{1}{\zeta^2} \left( \sum_{p=1}^{\infty} \frac{1}{p!} \left(\frac{1}{2}\zeta^2\right)^p \right) d\zeta \\
&= \frac{1}{2}x + \frac{1}{\sqrt{2\pi}} x \int_c^x \frac{1}{\zeta^2} \left( e^{\frac{1}{2}\zeta^2} - 1 \right) d\zeta \\
&= \frac{1}{2}x \left( \operatorname{erfi}\left(\frac{x}{\sqrt{2}}\right) + 1 - \operatorname{erfi}\left(\frac{c}{\sqrt{2}}\right) + \frac{1}{\sqrt{2\pi}} \frac{2}{c} \left( e^{\frac{1}{2}c^2} - 1 \right) \right) + \frac{1}{\sqrt{2\pi}} \left( 1 - e^{\frac{1}{2}x^2} \right)
\end{aligned}$$

Select  $c$  so that the first derivative is  $\frac{1}{2}x$ :

$$-\operatorname{erfi}\left(\frac{c}{\sqrt{2}}\right) + \frac{1}{\sqrt{2\pi}} \frac{2}{c} \left( e^{\frac{1}{2}c^2} - 1 \right) = 0 \text{ if } c = 0$$

Hence:

$$\begin{aligned}
\bar{s}^{[\text{ReLU}]}(x) &\triangleq \sum_{k=1}^{\infty} |a_k^{[\text{ReLU}]}| (1+x)^k - \sum_{k=1}^{\infty} |a_k^{[\text{ReLU}]}| \\
&= \frac{1}{2}(1+x) \left( \operatorname{erfi}\left(\frac{1+x}{\sqrt{2}}\right) + 1 \right) + \frac{1}{\sqrt{2\pi}} \left( 1 - e^{\frac{1}{2}(1+x)^2} \right) - \frac{1}{2} \left( \operatorname{erfi}\left(\frac{1}{\sqrt{2}}\right) + 1 \right) - \frac{1}{\sqrt{2\pi}} \left( 1 - e^{\frac{1}{2}} \right) \\
&= \frac{1}{2}x \left( \operatorname{erfi}\left(\frac{1+x}{\sqrt{2}}\right) + 1 \right) + \frac{1}{\sqrt{2\pi}} \left( e^{\frac{1}{2}} - e^{\frac{1}{2}(1+x)^2} \right) + \frac{1}{2} \left( \operatorname{erfi}\left(\frac{1+x}{\sqrt{2}}\right) - \operatorname{erfi}\left(\frac{1}{\sqrt{2}}\right) \right)
\end{aligned} \tag{28}$$

## B PROOF OF THE GLOBAL DUAL MODEL

Here we prove the validity of the global dual model presented in the paper. Recall that the dual model has the form (Theorem 1, equation (10)):

$$\mathbf{f}(\mathbf{x}; \Theta) = \langle \Psi(\Theta), \phi(\mathbf{x}) \rangle_{\mathbf{g}} \quad (29)$$

where, as per Figure 1,  $\Psi = \Psi^{[D-1]}$ ,  $\phi = \phi^{[D-1]}$ ,  $\mathbf{g} = \mathbf{g}^{[D-1]}$  and, given the base case  $\Psi^{[-1]}(\Theta) = \mathbf{1}_n$ ,  $\phi^{[-1]}(\mathbf{x}) = \mathbf{x}$ ,  $\mathbf{g}^{[-1]} = \mathbf{1}_n$ , the recursive definition of the feature maps and metric is proposed:

$$\begin{aligned} \tilde{\Psi}_{:i_{\tilde{j}}}^{[\tilde{j},j]}(\Theta) &= \tilde{\phi} \left[ a_{(0)k}^{[\tilde{j},j]} \right]^{\frac{1}{2}} \left[ \binom{k}{l}^{\frac{1}{2}} \left( \sqrt{\frac{1}{\tilde{s}^{[\tilde{j},j]} - 1} (\tilde{\phi}^2)} \Psi_{:i_{\tilde{j}}}^{[\tilde{j}]}(\Theta) \right)^{\otimes l} \right]_{1 \leq l \leq k} \Bigg]_{k \geq 1} \\ \tilde{\phi}^{[\tilde{j},j]}(\mathbf{x}) &= \frac{1}{\tilde{\phi}} \left[ a_{(0)k}^{[\tilde{j},j]} \right]^{\langle \frac{1}{2} \rangle} \left[ \binom{k}{l}^{\frac{1}{2}} \left( \sqrt{\frac{1}{\tilde{s}^{[\tilde{j},j]} - 1} (\tilde{\phi}^2)} \phi^{[\tilde{j}]}(\mathbf{x}) \right)^{\otimes l} \right]_{1 \leq l \leq k} \Bigg]_{k \geq 1} \\ \tilde{\mathbf{g}}^{[\tilde{j},j]} &= \left[ \text{He}_{k-l} \tilde{\mathbf{g}}^{[\tilde{j}]} \otimes l \right]_{1 \leq l \leq k} \Bigg]_{k \geq 1} \end{aligned} \quad (30)$$

$\forall \tilde{j} \in \tilde{\mathbb{P}}^{[j]}$  (the feature map transforms associated with the edges of the graph) and:

$$\begin{aligned} \Psi^{[j]}(\Theta) &= \sqrt{\gamma^2 + 1} \left[ \begin{array}{c} \mathbf{b}^{[j]\text{T}} + \mathbf{v}_{\tau}^{[j]\text{T}} \\ \text{diag} \left( \sqrt{\frac{\tilde{H}^{[j]}}{H^{[j]}}} \tilde{\Psi}^{[\tilde{j},j]}(\Theta) \mathbf{W}^{[\tilde{j},j]} \right)_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \end{array} \right] \\ \phi^{[j]}(\mathbf{x}) &= \sqrt{\frac{1}{\gamma^2 + 1}} \left[ \begin{array}{c} \gamma \\ \left[ \sqrt{\frac{\tilde{H}^{[j]}}{H^{[j]}}} \tilde{\phi}^{[\tilde{j},j]}(\mathbf{x}) \right]_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \end{array} \right] \\ \mathbf{g}^{[j]} &= \left[ \begin{array}{c} 1 \\ \left[ \tilde{\mathbf{g}}^{[\tilde{j},j]} \right]_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \end{array} \right] \end{aligned} \quad (31)$$

(the feature map transforms associated with the nodes of the graph) where  $\mathbf{v}_{\tau}^{[j]} = \sum_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \frac{\tau^{[\tilde{j},j]}(0)}{\gamma} \mathbf{W}^{[\tilde{j},j]\text{T}} \mathbf{1}_{H^{[\tilde{j}]}}$ . Our approach to demonstrating that this is true is to prove that, given some input  $\mathbf{x} \in \mathbb{X}$  then, for all edges ( $\tilde{j} \rightarrow j$ ):

$$\tilde{\mathbf{x}}^{[\tilde{j},j]} - \mathbf{1}_{H^{[\tilde{j}]}} \tau^{[\tilde{j},j]}(0) = \left\langle \tilde{\Psi}^{[\tilde{j},j]}(\Theta), \tilde{\phi}^{[\tilde{j},j]}(\mathbf{x}) \right\rangle_{\tilde{\mathbf{g}}^{[\tilde{j},j]}} \quad (32)$$

and likewise for all nodes  $j$ :

$$\mathbf{x}^{[j]} = \langle \Psi^{[j]}(\Theta), \phi^{[j]}(\mathbf{x}) \rangle_{\mathbf{g}^{[j]}} \quad (33)$$

**Base case:** By the definition of the base case, we have:

$$\mathbf{x}^{[-1]} = \langle \Psi^{[-1]}(\Theta), \phi^{[-1]}(\mathbf{x}) \rangle_{\mathbf{g}^{[-1]}} = \mathbf{x}$$

**Node case:** Assume (32) is true. Then, using (31), we have that:

$$\begin{aligned} \langle \Psi^{[j]}(\Theta), \phi^{[j]}(\mathbf{x}) \rangle_{\mathbf{g}^{[j]}} &= \gamma \mathbf{b}^{[j]} + \sum_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \left( \mathbf{W}^{[\tilde{j},j]\text{T}} \mathbf{1}_{H^{[\tilde{j}]}} \tau^{[\tilde{j},j]}(0) + \mathbf{W}^{[\tilde{j},j]\text{T}} \tilde{\mathbf{x}}^{[\tilde{j},j]} - \mathbf{W}^{[\tilde{j},j]\text{T}} \mathbf{1}_{H^{[\tilde{j}]}} \tau^{[\tilde{j},j]}(0) \right) \\ &= \gamma \mathbf{b}^{[j]} + \sum_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \mathbf{W}^{[\tilde{j},j]\text{T}} \tilde{\mathbf{x}}^{[\tilde{j},j]} \\ &= \mathbf{x}^{[j]} \end{aligned}$$

**Edge case:** Assume (33) is true. Then, using (30), we have that:

$$\begin{aligned}
\left\langle \tilde{\Psi}^{[\bar{j},j]}(\Theta), \tilde{\phi}^{[\bar{j},j]}(\mathbf{x}) \right\rangle_{\tilde{\mathbf{g}}^{[\bar{j},j]}} &= \left[ \left\langle \tilde{\Psi}_{:i_{\bar{j}}}^{[\bar{j},j]}(\Theta), \tilde{\phi}^{[\bar{j},j]}(\mathbf{x}) \right\rangle_{\tilde{\mathbf{g}}^{[\bar{j},j]}} \right]_{i_{\bar{j}}} \\
&= \left[ \sum_{k \geq 1} a_{(0)k}^{[\bar{j},j]} \sum_{1 \leq l \leq k} \binom{k}{l} \text{He}_{k-l} \left\langle \Psi_{:i_{\bar{j}}}^{[j]}(\Theta)^{\otimes l}, \phi^{[j]}(\mathbf{x})^{\otimes l} \right\rangle_{\mathbf{g}^{[j] \otimes l}} \right]_{i_{\bar{j}}} \\
&= \left[ \sum_{k \geq 1} a_{(0)k}^{[\bar{j},j]} \sum_{1 \leq l \leq k} \binom{k}{l} \text{He}_{k-l} \left\langle \Psi_{:i_{\bar{j}}}^{[j]}(\Theta), \phi^{[j]}(\mathbf{x}) \right\rangle_{\mathbf{g}^{[j]}}^l \right]_{i_{\bar{j}}} \\
&= \left[ \sum_{k \geq 1} a_{(0)k}^{[\bar{j},j]} \sum_{1 \leq l \leq k} \binom{k}{l} \text{He}_{k-l} x_{i_{\bar{j}}}^{[j]l} \right]_{i_{\bar{j}}} \\
&= \left[ \sum_{k \geq 0} a_{(0)k}^{[\bar{j},j]} \sum_{0 \leq l \leq k} \binom{k}{l} \text{He}_{k-l} x_{i_{\bar{j}}}^{[j]l} - \sum_{k \geq 0} a_{(0)k}^{[\bar{j},j]} \sum_{0 \leq l \leq k} \binom{k}{l} \text{He}_{k-l} 0^l \right]_{i_{\bar{j}}} \\
&= \left[ \tau^{[\bar{j},j]}(x_{i_{\bar{j}}}^{[j]}) - \tau^{[\bar{j},j]}(0) \right]_{i_{\bar{j}}} \\
&= \tilde{\mathbf{x}}^{[\bar{j},j]} - \mathbf{1}_{H^{[\bar{j}]}} \tau^{[\bar{j},j]}(0)
\end{aligned}$$

The desired result (29) then follows by identifying the output node  $j = D - 1$ .

## B.1 NORM-BOUNDS FOR THE GLOBAL DUAL MODEL

Our proof of the norm-bounds of the global model follows the same model as our proof of the validity of said model. We want to prove the bounds

$$\begin{aligned}
\left\| \tilde{\phi}^{[\bar{j},j]}(\mathbf{x}) \right\|_2^2 &\in \left[ \tilde{\phi}_{\downarrow}^{[\bar{j},j]2} = \frac{1}{\tilde{\phi}^2} \tilde{s}^{[\bar{j},j]} \left( \tilde{s}^{[\bar{j},j]-1} (\tilde{\phi}^2) \phi_{\downarrow}^{[j]2} \right), \tilde{\phi}^{[\bar{j},j]2} = 1 \right] \\
\left\| \tilde{\Psi}^{[\bar{j},j]}(\Theta) \right\|_2^2 &\leq \tilde{\psi}^{[\bar{j},j]2} = \tilde{\phi}^2 \tilde{s}^{[\bar{j},j]} \left( \frac{1}{\tilde{s}^{[\bar{j},j]-1} (\tilde{\phi}^2)} \psi^{[j]2} \right) \\
\left\| \tilde{\Psi}^{[\bar{j},j]}(\Theta) \right\|_{\text{He}[\tau]}^2 &\leq \tilde{\psi}^{[\bar{j},j]2} = \sup_{\substack{\phi_{\downarrow}^{[j]} \leq \phi^{[j]} \leq 1 \\ -\psi^{[j]} \leq \psi^{[j]} \leq \psi^{[j]}}} \left\{ \frac{\tilde{\phi}^2 \tau^{[\bar{j},j]} (\phi^{[j]} \psi^{[j]})^2}{\tilde{s}^{[\bar{j},j]} (\tilde{s}^{[\bar{j},j]-1} (\tilde{\phi}^2) \phi^{[j]2})} \right\}
\end{aligned} \tag{34}$$

$$\begin{aligned}
\left\| \phi^{[j]}(\mathbf{x}) \right\|_2^2 &\in \left[ \phi_{\downarrow}^{[j]2} = \frac{1}{\gamma^2 + 1} \left( \gamma^2 + \sum_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \frac{H^{[\tilde{j}]}}{H^{[j]}} \tilde{\phi}_{\downarrow}^{[\tilde{j},j]2} \right), \phi^{[j]2} = 1 \right] \\
\left\| \Psi^{[j]}(\Theta) \right\|_2^2 &\leq \psi^{[j]2} = (\gamma^2 + 1) \left( \left( \beta^{[j]} + \frac{\mu^{[j]} |\tau^{[\bar{j},j]}(0)|}{\gamma} \right)^2 + \sum_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \frac{\tilde{H}^{[\tilde{j}]}}{H^{[j]}} \tilde{\psi}^{[\tilde{j},j]2} \mu^{[\tilde{j},j]2} \right) \\
\left\| \Psi^{[j]}(\Theta) \right\|_{\text{He}[\tau]}^2 &\leq \psi^{[j]2} = (\gamma^2 + 1) \left( \left( \beta^{[j]} + \frac{\mu^{[j]} |\tau^{[\bar{j},j]}(0)|}{\gamma} \right)^2 + \sum_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \frac{\tilde{H}^{[\tilde{j}]}}{H^{[j]}} \tilde{\psi}^{[\tilde{j},j]2} \mu^{[\tilde{j},j]2} \right)
\end{aligned} \tag{35}$$

with the base-cases  $\|\phi^{[-1]}(\mathbf{x})\|_2^2 \in [\phi_{\downarrow}^{[-1]2} = 0, \phi^{[-1]2} = 1]$ ,  $\|\Psi^{[-1]}(\Theta)\|_2^2 \leq \psi^{[-1]2} = 1$  and  $\|\Psi^{[-1]}(\Theta)\|_{\text{He}[\tau]}^2 \leq \psi^{[-1]2} = 1$ . We proceed as follows:

**Base case:** By the definition of the base case, using our assumptions, we have:

$$\|\phi^{[-1]}(\mathbf{x})\|_2^2 = \|\mathbf{x}\|_2^2 \leq 1 = \phi^{[-1]2}$$

$$\|\Psi^{[-1]}(\Theta)\|_2^2 = \|\mathbf{I}_n\|_2^2 \leq 1 = \psi^{[-1]2}$$

$$\|\Psi^{[-1]}(\Theta)\|_{\text{He}[\tau]}^2 = \sup_{\mathbf{x} \in \mathbb{X}} \left\| \left\langle \frac{\phi^{[-1]}(\mathbf{x})}{\|\phi^{[-1]}(\mathbf{x})\|_2}, \Psi^{[-1]}(\Theta) \right\rangle_{\mathbf{g}} \right\|_2^2 = \sup_{\mathbf{x} \in \mathbb{X}} \left\| \frac{\phi^{[-1]}(\mathbf{x})}{\|\phi^{[-1]}(\mathbf{x})\|_2} \right\|_2^2 \leq 1 = \psi^{[-1]2}$$

**Node case:** Assume (34) is true. Then, using our assumptions, we have that:

$$\begin{aligned}
1242 \quad & \|\phi^{[j]}(\mathbf{x})\|_2^2 = \frac{1}{\gamma^2+1} \left\| \left[ \left[ \sqrt{\frac{H^{[j]}}{H^{[j]}}} \tilde{\Phi}^{[j]}(\mathbf{x}) \right]_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \right] \right\|_2^2 \\
1243 \quad & = \frac{1}{\gamma^2+1} \left( \gamma^2 + \sum_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \frac{H^{[j]}}{H^{[j]}} \|\tilde{\Phi}^{[j]}(\mathbf{x})\|_2^2 \right) \\
1244 \quad & \in \left[ \phi_{\downarrow}^{[j]2} = \frac{1}{\gamma^2+1} \left( \gamma^2 + \sum_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \frac{H^{[j]}}{H^{[j]}} \tilde{\phi}_{\downarrow}^{[j]2} \right), \phi^{[j]2} = 1 \right] \\
1245 \quad & \\
1246 \quad & \|\Psi^{[j]}(\Theta)\|_2^2 = (\gamma^2+1) \left\| \left[ \begin{array}{c} \mathbf{b}^{[j]T} + \mathbf{v}_{\tau}^{[j]T} \\ \text{diag} \left( \sqrt{\frac{\tilde{H}^{[j]}}{H^{[j]}}} \tilde{\Psi}^{[j]}(\Theta) \mathbf{W}^{[j]} \right) \end{array} \right] \right\|_2^2 \\
1247 \quad & \leq (\gamma^2+1) \left( \|\mathbf{b}^{[j]}\|_2 + \|\mathbf{v}_{\tau}^{[j]}\|_2 \right)^2 + \sum_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \frac{\tilde{H}^{[j]}}{H^{[j]}} \|\tilde{\Psi}^{[j]}(\Theta) \mathbf{W}^{[j]}\|_2^2 \\
1248 \quad & \leq (\gamma^2+1) \left( \left( \beta^{[j]} + \frac{\mu^{[j]} |\tau^{[j]}(0)|}{\gamma} \right)^2 + \sum_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \frac{\tilde{H}^{[j]}}{H^{[j]}} \tilde{\psi}^{[j]2} \mu^{[j]2} \right) = \psi^{[j]2} \\
1249 \quad & \\
1250 \quad & \|\Psi^{[j]}(\Theta)\|_{\text{He}[\tau]}^2 = (\gamma^2+1) \left\| \left[ \begin{array}{c} \mathbf{b}^{[j]T} + \mathbf{v}_{\tau}^{[j]T} \\ \text{diag} \left( \sqrt{\frac{\tilde{H}^{[j]}}{H^{[j]}}} \tilde{\Psi}^{[j]}(\Theta) \mathbf{W}^{[j]} \right) \end{array} \right] \right\|_2^2 \\
1251 \quad & \leq (\gamma^2+1) \left( \|\mathbf{b}^{[j]}\|_2 + \|\mathbf{v}_{\tau}^{[j]}\|_2 \right)^2 + \sum_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \frac{\tilde{H}^{[j]}}{H^{[j]}} \|\tilde{\Psi}^{[j]}(\Theta) \mathbf{W}^{[j]}\|_{\text{He}[\tau]}^2 \\
1252 \quad & \leq (\gamma^2+1) \left( \left( \beta^{[j]} + \frac{\mu^{[j]} |\tau^{[j]}(0)|}{\gamma} \right)^2 + \sum_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \frac{\tilde{H}^{[j]}}{H^{[j]}} \tilde{\psi}^{[j]2} \mu^{[j]2} \right) = \tilde{\psi}^{[j]2} \\
1253 \quad & \\
1254 \quad & \\
1255 \quad & \\
1256 \quad & \\
1257 \quad & \\
1258 \quad & \\
1259 \quad & \\
1260 \quad & \\
1261 \quad & \\
1262 \quad & \\
1263 \quad & \\
1264 \quad & \\
1265 \quad & \\
1266 \quad & \\
1267 \quad & \\
1268 \quad & \\
1269 \quad & \\
1270 \quad &
\end{aligned}$$

**Edge case:** Assume (35) is true. Then, using our assumptions and the definition (12), and using that the magnitude function is increasing on  $\mathbb{R}_+$ , we have that:

$$\begin{aligned}
1271 \quad & \|\tilde{\Phi}^{[j]}(\mathbf{x})\|_2^2 = \frac{1}{\tilde{\phi}^2} \left\| \left[ \left[ a_{(0)k}^{[j]} \left\langle \frac{1}{2} \right\rangle \left[ \binom{k}{l} \frac{1}{2} \left( \sqrt{\tilde{s}^{[j]}-1} \left( \tilde{\phi}^2 \right) \phi^{[j]}(\mathbf{x}) \right)^{\otimes l} \right]_{1 \leq l \leq k} \right]_{k \geq 1} \right] \right\|_2^2 \\
1272 \quad & = \frac{1}{\tilde{\phi}^2} \sum_{k \geq 1} \left| a_{(0)k}^{[j]} \right| \sum_{1 \leq l \leq k} \binom{k}{l} \left( \tilde{s}^{[j]}-1 \right) \left( \tilde{\phi}^2 \right) \|\phi^{[j]}(\mathbf{x})\|_2^2 \\
1273 \quad & = \frac{1}{\tilde{\phi}^2} \tilde{s}^{[j]} \left( \tilde{s}^{[j]}-1 \right) \left( \tilde{\phi}^2 \right) \|\phi^{[j]}(\mathbf{x})\|_2^2 \\
1274 \quad & \in \left[ \frac{1}{\tilde{\phi}^2} \tilde{s}^{[j]} \left( \tilde{s}^{[j]}-1 \right) \left( \tilde{\phi}^2 \right) \phi_{\downarrow}^{[j]2}, \frac{1}{\tilde{\phi}^2} \tilde{s}^{[j]} \left( \tilde{s}^{[j]}-1 \right) \left( \tilde{\phi}^2 \right) \phi^{[j]2} \right] \\
1275 \quad & \in \left[ \frac{1}{\tilde{\phi}^2} \tilde{s}^{[j]} \left( \tilde{s}^{[j]}-1 \right) \left( \tilde{\phi}^2 \right) \phi_{\downarrow}^{[j]2} = \tilde{\phi}_{\downarrow}^{[j]2}, \tilde{\phi}^{[j]2} = 1 \right] \\
1276 \quad & \\
1277 \quad & \|\tilde{\Psi}^{[j]}(\Theta)\|_2^2 = \tilde{\phi}^2 \max_{i_{\tilde{j}}} \left\| \left[ \left[ \left| a_{(0)k}^{[j]} \right|^{\frac{1}{2}} \left[ \binom{k}{l} \frac{1}{2} \left( \sqrt{\frac{1}{\tilde{s}^{[j]}-1} \left( \tilde{\phi}^2 \right)} \Psi^{[j]}(\Theta) \right)^{\otimes l} \right]_{1 \leq l \leq k} \right]_{k \geq 1} \right] \right\|_2^2 \\
1278 \quad & = \tilde{\phi}^2 \max_{i_{\tilde{j}}} \tilde{s}^{[j]} \left( \frac{1}{\tilde{s}^{[j]}-1} \left( \tilde{\phi}^2 \right) \|\Psi^{[j]}(\Theta)\|_2^2 \right) \\
1279 \quad & = \tilde{\phi}^2 \tilde{s}^{[j]} \left( \frac{1}{\tilde{s}^{[j]}-1} \left( \tilde{\phi}^2 \right) \|\Psi^{[j]}(\Theta)\|_2^2 \right) \\
1280 \quad & \leq \tilde{\phi}^2 \tilde{s}^{[j]} \left( \frac{1}{\tilde{s}^{[j]}-1} \left( \tilde{\phi}^2 \right) \psi^{[j]2} \right) \\
1281 \quad & \\
1282 \quad & \\
1283 \quad & \\
1284 \quad & \\
1285 \quad & \\
1286 \quad & \\
1287 \quad & \\
1288 \quad & \\
1289 \quad & \\
1290 \quad & \\
1291 \quad & \\
1292 \quad & \\
1293 \quad & \\
1294 \quad & \\
1295 \quad & \|\tilde{\Psi}^{[j]}(\Theta)\|_{\text{He}[\tau]}^2 \leq \tilde{\psi}^{[j]2} = \sup_{\substack{\phi_{\downarrow}^{[j]} \leq \phi^{[j]} \leq 1 \\ -\psi^{[j]} \leq \psi^{[j]} \leq \psi^{[j]}}} \left\{ \frac{\tilde{\phi}^2 \tau^{[j]} \left( \phi^{[j]} \psi^{[j]} \right)^2}{\tilde{s}^{[j]} \left( \tilde{s}^{[j]}-1 \right) \left( \tilde{\phi}^2 \right) \phi^{[j]2}} \right\}
\end{aligned}$$

where the final bound in this sequence is simply the definition of the norm in question with the range of the supremum expanded to the known bound on this range.

The desired result follows by identifying the output node  $j = D - 1$ . We note that the bounds  $\tilde{\phi} \in \mathbb{R}_+$  may be chosen arbitrarily here.

## C PROOF OF THE LOCAL DUAL MODEL

We now repeat the proof from the previous section B for the local model. Recall that the dual model has the form (Theorem 1, equation (10)):

$$\begin{aligned} \Delta \mathbf{f}(\mathbf{x}; \Delta \Theta) &= \langle \Psi_{\Delta}(\Delta \Theta), \phi_{\Delta}(\mathbf{x}) \rangle_{\mathbf{G}_{\Delta}(\mathbf{x})} \\ &= \left[ \langle \Psi_{\Delta:i_{D-1}}(\Delta \Theta), \phi_{\Delta}(\mathbf{x}) \rangle_{\mathbf{G}_{\Delta:i_{D-1}}(\mathbf{x})} \right]_{i_{D-1}} \end{aligned} \quad (36)$$

where, as per Figure 2,  $\Psi_{\Delta} = \Psi_{\Delta}^{[D-1]}$ ,  $\phi_{\Delta} = \phi_{\Delta}^{[D-1]}$ ,  $\mathbf{G}_{\Delta} = \mathbf{G}_{\Delta}^{[D-1]}$  and, given the base case  $\Psi_{\Delta}^{[-1]}(\Delta \Theta) = \mathbf{0}_{0 \times n}$ ,  $\phi_{\Delta}^{[-1]}(\mathbf{x}) = \mathbf{0}_0$ ,  $\mathbf{G}_{\Delta}^{[-1]}(\mathbf{x}) = \mathbf{1}_{0 \times n}$  the recursive definition of the feature maps is proposed as:

$$\begin{aligned} \tilde{\phi}_{\Delta}^{[\tilde{j},j]}(\mathbf{x}) &= \left[ \frac{1}{\eta^k} \left[ \left( \sqrt{\frac{\rho_{(\tilde{\omega})\eta}^{[\tilde{j},j]^2}}{\gamma^2+1}} \phi_{\Delta}^{[\tilde{j},j]}(\mathbf{x}) \right)^{\otimes l} \right]_{1 \leq l \leq k} \right]_{k \geq 1} \\ \tilde{\Psi}_{\Delta:i_{\tilde{j}}}^{[\tilde{j},j]}(\Theta) &= \left[ \eta^k \left[ \left( \sqrt{\frac{\gamma^2+1}{\rho_{(\tilde{\omega})\eta}^{[\tilde{j},j]^2}}} \Psi_{\Delta:i_{\tilde{j}}}^{[\tilde{j},j]}(\Theta) \right)^{\otimes l} \right]_{1 \leq l \leq k} \right]_{k \geq 1} \\ \tilde{\mathbf{G}}_{\Delta:i_{\tilde{j}}}^{[\tilde{j},j]}(\mathbf{x}) &= \left[ a_{\binom{[\tilde{j},j]}{x_{i_{\tilde{j}}}}^{[\tilde{j},j]}} \left[ \binom{k}{l} \text{He}_{k-l} \mathbf{G}_{\Delta:i_{\tilde{j}}}^{[\tilde{j},j]}(\mathbf{x})^{\otimes l} \right]_{1 \leq l \leq k} \right]_{k \geq 1} \end{aligned} \quad (37)$$

$\forall \tilde{j} \in \tilde{\mathbb{P}}^{[j]}$  (the feature map transforms associated with the edges of the graph) and:

$$\begin{aligned} \phi_{\Delta}^{[j]}(\mathbf{x}) &= \left[ \begin{array}{c} \frac{1}{\sqrt{2\tilde{p}^{[j]}}} \left[ \frac{\gamma}{\tilde{\omega}^{[\tilde{j},j]}} \tilde{\mathbf{x}}^{[\tilde{j},j]} \right]_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \\ \frac{1}{\sqrt{2\tilde{p}^{[j]}}} \left[ \frac{1}{\tilde{\psi}^{[\tilde{j},j]}} \tilde{\phi}_{\Delta}^{[\tilde{j},j]}(\mathbf{x}) \right]_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \\ \Delta \mathbf{b}^{[j]\text{T}} \end{array} \right] \\ \Psi_{\Delta}^{[j]}(\Theta) &= \left[ \begin{array}{c} \sqrt{2\tilde{p}^{[j]}} \text{diag} \left( \tilde{\omega}^{[\tilde{j},j]} \mathbf{I}_{H^{[\tilde{j}]}} \right) \Delta \mathbf{W}^{[j]} \\ \sqrt{2\tilde{p}^{[j]}} \text{diag} \left( \tilde{\psi}^{[\tilde{j},j]} \tilde{\Psi}_{\Delta}^{[\tilde{j},j]}(\Theta) \right) (\mathbf{W}^{[j]} + \Delta \mathbf{W}^{[j]}) \end{array} \right]_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \\ \mathbf{G}_{\Delta}^{[j]}(\mathbf{x}) &= \left[ \begin{array}{c} \mathbf{1}_{H^{[\tilde{j}]}}^{\text{T}} \\ \text{diag} \left( \mathbf{I}_{H^{[\tilde{j}]}} \mathbf{1}_{\tilde{H}^{[j]}} \mathbf{1}_{H^{[\tilde{j}]}}^{\text{T}} \right)_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \\ \text{diag} \left( \tilde{\mathbf{G}}_{\Delta}^{[\tilde{j},j]}(\mathbf{x}) \right) \mathbf{1}_{\tilde{H}^{[j]}} \mathbf{1}_{H^{[\tilde{j}]}}^{\text{T}} \end{array} \right]_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \end{aligned} \quad (38)$$

(the feature map transforms associated with the nodes of the graph). Our approach to demonstrating that this is true is to prove that, given some input  $\mathbf{x} \in \mathbb{X}$  then, for all edges ( $\tilde{j} \rightarrow j$ ):

$$\Delta \tilde{\mathbf{x}}^{[\tilde{j},j]} = \left\langle \tilde{\Psi}_{\Delta}^{[\tilde{j},j]}(\Delta \Theta), \tilde{\phi}_{\Delta}^{[\tilde{j},j]}(\mathbf{x}) \right\rangle_{\tilde{\mathbf{G}}_{\Delta}^{[\tilde{j},j]}(\mathbf{x})} \quad (39)$$

and likewise for all nodes  $j$ :

$$\Delta \mathbf{x}^{[j]} = \left\langle \Psi_{\Delta}^{[j]}(\Delta \Theta), \phi_{\Delta}^{[j]}(\mathbf{x}) \right\rangle_{\mathbf{G}_{\Delta}^{[j]}(\mathbf{x})} \quad (40)$$

1350 **Base case:** By the definition of the base case, we have:

$$1351 \Delta \mathbf{x}^{[-1]} = \left\langle \Psi_{\Delta}^{[-1]}(\Delta\Theta), \phi_{\Delta}^{[-1]}(\mathbf{x}) \right\rangle_{\mathbf{G}_{\Delta}^{[-1]}(\mathbf{x})} = \mathbf{0}_n$$

1352  
1353  
1354 **Node case:** Assume (39) is true. Then, using (38), we have that:

$$1355 \left\langle \Psi_{\Delta}^{[j]}(\Delta\Theta), \phi_{\Delta}^{[j]}(\mathbf{x}) \right\rangle_{\mathbf{G}_{\Delta}^{[j]}(\mathbf{x})} = \gamma \Delta \mathbf{b}^{[j]} + \sum_{\tilde{j} \in \mathbb{P}^{[j]}} \left( \Delta \mathbf{W}^{[\tilde{j},j]} \mathbf{T} \tilde{\mathbf{x}}^{[\tilde{j},j]} + \mathbf{W}^{[\tilde{j},j]} \mathbf{T} \Delta \tilde{\mathbf{x}}^{[\tilde{j},j]} + \Delta \mathbf{W}^{[\tilde{j},j]} \mathbf{T} \Delta \tilde{\mathbf{x}}^{[\tilde{j},j]} \right)$$

$$1356 = \gamma (\mathbf{b}^{[j]} + \Delta \mathbf{b}^{[j]}) + \sum_{\tilde{j} \in \mathbb{P}^{[j]}} (\mathbf{W}^{[\tilde{j},j]} + \Delta \mathbf{W}^{[\tilde{j},j]})^{\mathbf{T}} (\tilde{\mathbf{x}}^{[\tilde{j},j]} + \Delta \tilde{\mathbf{x}}^{[\tilde{j},j]}) - \gamma \mathbf{b}^{[j]} - \sum_{\tilde{j} \in \mathbb{P}^{[j]}} \mathbf{W}^{[\tilde{j},j]} \mathbf{T} \tilde{\mathbf{x}}^{[\tilde{j},j]}$$

$$1357 = (\mathbf{x}^{[j]} + \Delta \mathbf{x}^{[j]}) - \mathbf{x}^{[j]} = \Delta \mathbf{x}^{[j]}$$

1358  
1359  
1360 **Edge case:** Assume (40) is true. Then, using (37) and (27), we have that:

$$1361 \left\langle \tilde{\Psi}_{\Delta}^{[\tilde{j},j]}(\Delta\Theta), \tilde{\phi}_{\Delta}^{[\tilde{j},j]}(\mathbf{x}) \right\rangle_{\tilde{\mathbf{G}}_{\Delta}^{[\tilde{j},j]}(\mathbf{x})} = \left[ \sum_{k \geq 1} a^{[\tilde{j},j]}(x_{i_{\tilde{j}}}^{[\tilde{j}]})^k \sum_{1 \leq l \leq k} \binom{k}{l} \text{He}_{k-l} \left\langle \Psi_{\Delta:i_{\tilde{j}}}^{[\tilde{j}]}(\Delta\Theta), \phi_{\Delta}^{[\tilde{j}]}(\mathbf{x}) \right\rangle_{\tilde{\mathbf{G}}_{\Delta}^{[\tilde{j},j]}(\mathbf{x})} \right]_{i_{\tilde{j}}}$$

$$1362 = \left[ \sum_{k \geq 1} a^{[\tilde{j},j]}(x_{i_{\tilde{j}}}^{[\tilde{j}]})^k \sum_{1 \leq l \leq k} \binom{k}{l} \text{He}_{k-l} \Delta x_{i_{\tilde{j}}}^{[j]l} \right]_{i_{\tilde{j}}}$$

$$1363 = \left[ \bar{\tau}^{[\tilde{j},j]}(\Delta x_{i_{\tilde{j}}}^{[j]}, x_{i_{\tilde{j}}}^{[\tilde{j}]}) \right]_{i_{\tilde{j}}}$$

$$1364 = \left[ \tau^{[\tilde{j},j]}(x_{i_{\tilde{j}}}^{[\tilde{j}]} + \Delta x_{i_{\tilde{j}}}^{[j]}) - \tau^{[\tilde{j},j]}(x_{i_{\tilde{j}}}^{[\tilde{j}]}) \right]_{i_{\tilde{j}}}$$

$$1365 = \Delta \tilde{\mathbf{x}}^{[\tilde{j},j]}$$

1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374 The desired result (36) then follows by identifying the output node  $j = D - 1$ .

### 1375 C.1 NORM-BOUNDS FOR THE LOCAL DUAL MODEL

1376  
1377 Our proof of the norm-bounds of the local model follows the same model as our proof of the validity  
1378 of said model. We want to prove the bounds

$$1379 \left\| \tilde{\phi}_{\Delta}^{[\tilde{j},j]}(\mathbf{x}) \odot \tilde{\mathbf{G}}_{\Delta:i_{\tilde{j}}}^{[\tilde{j},j]}(\mathbf{x}) \right\|_2^2 \leq \tilde{\phi}_{\Delta}^{[\tilde{j},j]2} = T_{(\tilde{\omega}^{[\tilde{j},j]})_{\eta}}^{[\tilde{j},j]2} \forall i_{\tilde{j}}$$

$$1380 \left\| \tilde{\Psi}_{\Delta}^{[\tilde{j},j]}(\Delta\Theta) \right\|_2^2 \leq \tilde{\psi}_{\Delta}^{[\tilde{j},j]2} = \hat{s}_{\eta} \left( \frac{\gamma^2 + 1}{\rho_{(\tilde{\omega}^{[\tilde{j},j]})_{\eta}}^{[\tilde{j},j]2}} \psi_{\Delta}^{[\tilde{j}]2} \right)$$

$$1381 \left\| \phi_{\Delta}^{[j]}(\mathbf{x}) \odot \mathbf{G}_{\Delta:i_j}^{[j]}(\mathbf{x}) \right\|_2^2 \leq \phi_{\Delta}^{[j]2} = \gamma^2 + 1 \forall i_j$$

$$1382 \left\| \Psi_{\Delta}^{[j]}(\Theta) \right\|_2^2 \leq \psi_{\Delta}^{[j]2} = \beta_{\Delta}^{[j]2} + 2\tilde{p}^{[j]} \left( \mu_{\Delta}^{[j]2} \tilde{\omega}^{[j]2} + (\mu^{[j]2} + \mu_{\Delta}^{[j]2}) \tilde{\psi}_{\Delta}^{[j]2} \right)$$

1383  
1384  
1385  
1386  
1387  
1388  
1389 with the base-cases  $\|\phi_{\Delta}^{[-1]}(\mathbf{x}) \odot \mathbf{G}_{\Delta:i_{-1}}^{[-1]}(\mathbf{x})\|_2^2 \leq \phi_{\Delta}^{[-1]2} = 0$  and  $\|\Psi_{\Delta}^{[-1]}(\Delta\Theta)\|_2^2 \leq \psi_{\Delta}^{[-1]2} = 0$ . We  
1390 proceed as follows:

1391 **Base case:** By the definition of the base case, using our assumptions, we have:

$$1392 \left\| \phi_{\Delta}^{[-1]}(\mathbf{x}) \odot \mathbf{G}_{\Delta:i_{-1}}^{[-1]}(\mathbf{x}) \right\|_2^2 = 0 = \phi^{[-1]2} \forall i_{-1}$$

$$1393 \left\| \Psi_{\Delta}^{[-1]}(\Delta\Theta) \right\|_2^2 = 0 = \psi^{[-1]2}$$

1394  
1395  
1396  
1397 **Node case:** Assume (41) is true. Then, using our assumptions, we have that:

$$1400 \left\| \phi_{\Delta}^{[j]}(\mathbf{x}) \odot \mathbf{G}_{\Delta:i_j}^{[j]}(\mathbf{x}) \right\|_2^2 = \gamma^2 + \frac{1}{2\tilde{p}^{[j]}} \sum_{\tilde{j} \in \mathbb{P}^{[j]}} \frac{1}{\tilde{\omega}^{[\tilde{j},j]}} \|\mathbf{x}^{[\tilde{j},j]}\|_2^2 + \frac{1}{2\tilde{p}^{[j]}} \sum_{\tilde{j} \in \mathbb{P}^{[j]}} \frac{1}{\tilde{\psi}_{\Delta}^{[\tilde{j},j]}} \left\| \tilde{\phi}_{\Delta}^{[\tilde{j},j]}(\mathbf{x}) \odot \tilde{\mathbf{G}}_{\Delta:i_{\tilde{j}}}^{[\tilde{j},j]}(\mathbf{x}) \right\|_2^2$$

$$1401 \leq \gamma^2 + 1 = \phi_{\Delta}^{[j]2}$$

$$\begin{aligned}
1404 \quad \|\Psi_{\Delta}^{[j]}(\Delta\Theta)\|_2^2 &= \|\Delta\mathbf{b}^{[j]T}\|_2^2 + 2\tilde{p}^{[j]} \max_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \tilde{\omega}^{[\tilde{j},j]2} \|\Delta\mathbf{W}^{[j]}\|_2^2 + 2\tilde{p}^{[j]} \max_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \tilde{\psi}_{\Delta}^{[\tilde{j},j]2} \|\mathbf{W}^{[j]} + \Delta\mathbf{W}^{[j]}\|_2^2 \\
1405 &\leq \beta_{\Delta}^{[j]2} + 2\tilde{p}^{[j]} \max_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \tilde{\omega}^{[\tilde{j},j]2} \mu_{\Delta}^{[j]2} + 2\tilde{p}^{[j]} \max_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \tilde{\psi}_{\Delta}^{[\tilde{j},j]2} (\mu^{[j]2} + \mu_{\Delta}^{[j]2}) \\
1406 &= \beta_{\Delta}^{[j]2} + 2\tilde{p}^{[j]} \tilde{\omega}^{[j]2} \mu_{\Delta}^{[j]2} + 2\tilde{p}^{[j]} \tilde{\psi}_{\Delta}^{[j]2} (\mu^{[j]2} + \mu_{\Delta}^{[j]2}) \\
1407 &= \beta_{\Delta}^{[j]2} + 2\tilde{p}^{[j]} (\mu_{\Delta}^{[j]2} \tilde{\omega}^{[j]2} + (\mu^{[j]2} + \mu_{\Delta}^{[j]2}) \tilde{\psi}_{\Delta}^{[j]2}) = \psi_{\Delta}^{[j]2}
\end{aligned}$$

**Edge case:** Assume (42) is true. Then, using our assumptions and the definition (21) of the rectified activation function and it's increasing (on  $\mathbb{R}_+$ ) envelope, we have that:

$$\begin{aligned}
1414 \quad \|\tilde{\phi}_{\Delta}^{[\tilde{j},j]}(\mathbf{x}) \odot \tilde{\mathbf{G}}_{\Delta:i_{\tilde{j}}}^{[\tilde{j},j]}(\mathbf{x})\|_2^2 &= \left\| \left[ \frac{a^{[\tilde{j},j]2}}{\eta^k} \left( \frac{x_{i_{\tilde{j}}}^{[\tilde{j}]}}{\eta^k} \right)^k \left( \binom{k}{l} \text{He}_{k-l} \left( \frac{\rho^{[\tilde{j},j]2}}{\tilde{\omega}^{[\tilde{j},j]2}} \right)^{\frac{l}{2}} \phi_{\Delta}^{[j]}(\mathbf{x})^{\otimes l} \odot \mathbf{G}_{\Delta:i_{\tilde{j}}}^{[j]}(\mathbf{x})^{\otimes l} \right) \right]_{1 \leq l \leq k} \right\|_{k \geq 1}^2 \\
1415 &= \left\| \left[ \frac{a^{[\tilde{j},j]2}}{\eta^k} \left( \frac{x_{i_{\tilde{j}}}^{[\tilde{j}]}}{\eta^k} \right)^k \left( \binom{k}{l} \text{He}_{k-l} \left( \left( \frac{\rho^{[\tilde{j},j]2}}{\tilde{\omega}^{[\tilde{j},j]2}} \right)^{\frac{l}{2}} \phi_{\Delta}^{[j]}(\mathbf{x}) \odot \mathbf{G}_{\Delta:i_{\tilde{j}}}^{[j]}(\mathbf{x}) \right)^{\otimes l} \right) \right]_{1 \leq l \leq k} \right\|_{k \geq 1}^2 \\
1416 &= \sum_{k \geq 1} \frac{a^{[\tilde{j},j]2}}{\eta^{2k}} \sum_{1 \leq l \leq k} \binom{k}{l}^2 \text{He}_{k-l}^2 \left( \frac{\rho^{[\tilde{j},j]2}}{\tilde{\omega}^{[\tilde{j},j]2}} \right) \left\| \phi_{\Delta}^{[j]}(\mathbf{x}) \odot \mathbf{G}_{\Delta:i_{\tilde{j}}}^{[j]}(\mathbf{x}) \right\|_2^2 \Big)^l \\
1417 &= \hat{\tau}_{\eta}^{[\tilde{j},j]} \left( \frac{\rho^{[\tilde{j},j]2}}{\tilde{\omega}^{[\tilde{j},j]2}} \left\| \phi_{\Delta}^{[j]}(\mathbf{x}) \odot \mathbf{G}_{\Delta:i_{\tilde{j}}}^{[j]}(\mathbf{x}) \right\|_2^2; x_{i_{\tilde{j}}}^{[\tilde{j}]}, x_{i_{\tilde{j}}}^{[j]} \right) \\
1418 &\leq \hat{\tau}_{\eta}^{[\tilde{j},j]} \left( \frac{\rho^{[\tilde{j},j]2}}{\tilde{\omega}^{[\tilde{j},j]2}} \left\| \phi_{\Delta}^{[j]}(\mathbf{x}) \odot \mathbf{G}_{\Delta:i_{\tilde{j}}}^{[j]}(\mathbf{x}) \right\|_2^2; \tilde{\omega}^{[\tilde{j},j]} \right) \\
1419 &\leq \hat{\tau}_{\eta}^{[\tilde{j},j]} \left( \rho_{(\tilde{\omega}^{[\tilde{j},j]})_{\eta}}^{[\tilde{j},j]2}; \tilde{\omega}^{[\tilde{j},j]} \right) \\
1420 &\leq T_{(\tilde{\omega}^{[\tilde{j},j]})_{\eta}}^{[\tilde{j},j]2} = \tilde{\phi}_{\Delta}^{[\tilde{j},j]2} \\
1421 & \\
1422 \quad \|\tilde{\Psi}_{\Delta}^{[\tilde{j},j]}(\Delta\Theta)\|_2^2 &= \max_{i_{\tilde{j}}} \left\| \left[ \eta^k \left[ \left( \frac{\gamma^2+1}{\rho_{(\tilde{\omega}^{[\tilde{j},j]})_{\eta}}^{[\tilde{j},j]2}} \right)^{\frac{l}{2}} (\Psi_{\Delta:i_{\tilde{j}}}^{[j]}(\Delta\Theta))^{\otimes l} \right]_{1 \leq l \leq k} \right] \right\|_{k \geq 1}^2 \\
1423 &= \max_{i_{\tilde{j}}} \hat{s}_{\eta} \left( \frac{\gamma^2+1}{\rho_{(\tilde{\omega}^{[\tilde{j},j]})_{\eta}}^{[\tilde{j},j]2}} \left\| \Psi_{\Delta:i_{\tilde{j}}}^{[j]}(\Delta\Theta) \right\|_2^2 \right) \\
1424 &= \hat{s}_{\eta} \left( \frac{\gamma^2+1}{\rho_{(\tilde{\omega}^{[\tilde{j},j]})_{\eta}}^{[\tilde{j},j]2}} \left\| \Psi_{\Delta}^{[j]}(\Delta\Theta) \right\|_2^2 \right) \\
1425 &\leq \hat{s}_{\eta} \left( \frac{\gamma^2+1}{\rho_{(\tilde{\omega}^{[\tilde{j},j]})_{\eta}}^{[\tilde{j},j]2}} \psi_{\Delta}^{[j]2} \right) = \tilde{\psi}_{\Delta}^{[\tilde{j},j]2}
\end{aligned} \tag{43}$$

The desired result follows by identifying the output node  $j = D - 1$ .

## C.2 CONVERGENCE REGION OF LOCAL DUAL MODEL

Here we prove the bounds on the change in weights and biases for which the local dual model holds as state in theorem 6. The goal is to place bounds on  $\mu_{\Delta}^{[j]}, \beta_{\Delta}^{[j]}$  to ensure that:

$$\frac{1}{\rho_{(\tilde{\omega}^{[\tilde{j},j]})_{\eta}}^{[\tilde{j},j]2}} \left\| \Psi_{\Delta}^{[j]}(\Delta\Theta) \right\|_2^2 \leq \frac{1}{\rho_{(\tilde{\omega}^{[\tilde{j},j]})_{\eta}}^{[\tilde{j},j]2}} \psi_{\Delta}^{[j]2} \leq R_{\eta}^2$$

$\forall \tilde{j} \in \tilde{\mathbb{P}}^{[j]}, \forall j$  in (43). This result suffices to ensure that  $\|\tilde{\Psi}_{\Delta}^{[\tilde{j},j]}(\Delta\Theta) \odot \tilde{\mathbf{G}}_{\Delta:i_{\tilde{j}}}^{[\tilde{j},j]}\|_2^2 \leq 1$  converges, and subsequently that both feature maps are convergent, allowing us to conclude that the model is well-defined (convergent) and enabling e.g. our bound on Rademacher complexity to be derived.

We aim to find the largest possible  $\mu_{\Delta}^{[\tilde{j},j]}, \beta_{\Delta}^{[j]}$  such that:

$$\frac{1}{\rho_{(\tilde{\omega}^{[\tilde{j},j]})_{\eta}}^{[\tilde{j},j]2}} \psi_{\Delta}^{[\tilde{j}]2} \leq R_{\eta}^2$$

$\forall j : \tilde{j} \in \tilde{\mathbb{P}}^{[j]}$ . So we require that:

$$\psi_{\Delta}^{[\tilde{j}]2} \leq \min_{j:\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \rho_{(\tilde{\omega}^{[\tilde{j},j]})_{\eta}}^{[\tilde{j},j]2} R_{\eta}^2 \quad (44)$$

Recall that:

$$\psi_{\Delta}^{[\tilde{j}]2} = (\gamma^2 + 1) \left( \beta_{\Delta}^{[\tilde{j}]2} + 2 \sum_{q \in \mathbb{P}^{[\tilde{j}]}} \frac{\tilde{H}^{[\tilde{j}]}}{H^{[q]}} \left( \mu_{\Delta}^{[q,\tilde{j}]2} \tilde{\omega}^{[q,\tilde{j}]2} + 2\mu^{[q,\tilde{j}]2} \tilde{\psi}_{\Delta}^{[q,\tilde{j}]2} \right) \right)$$

Thus our condition becomes:

$$\left( \beta_{\Delta}^{[\tilde{j}]2} + 2 \sum_{q \in \mathbb{P}^{[\tilde{j}]}} \frac{\tilde{H}^{[\tilde{j}]}}{H^{[q]}} \mu_{\Delta}^{[q,\tilde{j}]2} \tilde{\omega}^{[q,\tilde{j}]2} \right) + 4 \sum_{q \in \mathbb{P}^{[\tilde{j}]}} \frac{\tilde{H}^{[\tilde{j}]}}{H^{[q]}} \mu^{[q,\tilde{j}]2} \tilde{\psi}_{\Delta}^{[q,\tilde{j}]2} \leq \min_{j:\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \frac{\rho_{(\tilde{\omega}^{[\tilde{j},j]})_{\eta}}^{[\tilde{j},j]2}}{\gamma^2 + 1} R_{\eta}^2$$

Selecting  $d^{[\tilde{j}]} \in (0, 1) \forall j$  we can simplify the requirement to:

$$(\gamma^2 + 1) \left( \beta_{\Delta}^{[\tilde{j}]2} + 2p^{[\tilde{j}]} \mu_{\Delta}^{[\tilde{j}]2} \tilde{\omega}^{[\tilde{j}]2} \right) \leq d^{[\tilde{j}]} \psi_{\Delta}^{[\tilde{j}]2}$$

and:

$$4p^{[\tilde{j}]} \mu^{[\tilde{j}]2} \tilde{\psi}_{\Delta}^{[\tilde{j}]2} \leq (1 - d^{[\tilde{j}]}) \min_{j:\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \frac{\rho_{(\tilde{\omega}^{[\tilde{j},j]})_{\eta}}^{[\tilde{j},j]2}}{\gamma^2 + 1} R_{\eta}^2$$

Re-ordering, we find:

$$4\tilde{p}^{[\tilde{j}]} \mu^{[\tilde{j}]2} \hat{s}_{\eta} \left( \frac{1}{\rho_{(\tilde{\omega}^{[\tilde{j},j]})_{\eta}}^{[\tilde{j},j]2}} \psi_{\Delta}^{[\tilde{j}]2} \right) \leq (1 - d^{[\tilde{j}]}) \psi_{\Delta}^{[\tilde{j}]2} \quad \forall \tilde{j} \in \tilde{\mathbb{P}}^{[j]}$$

and after cleaning up and re-indexing:

$$\psi_{\Delta}^{[j]2} \leq \rho_{(\tilde{\omega}^{[j,\tilde{j}])_{\eta}}^{[j,\tilde{j}]2}} \hat{s}_{\eta}^{-1} \left( \frac{1}{4\tilde{p}^{[\tilde{j}]} \mu^{[\tilde{j}]2}} (1 - d^{[\tilde{j}]}) \psi_{\Delta}^{[\tilde{j}]2} \right) \quad \forall \tilde{j} : j \in \tilde{\mathbb{P}}^{[\tilde{j}]}$$

so, overall:

$$\psi_{\Delta}^{[j]2} \leq \min_{\tilde{j}:j \in \tilde{\mathbb{P}}^{[\tilde{j}]}} \rho_{(\tilde{\omega}^{[j,\tilde{j}])_{\eta}}^{[j,\tilde{j}]2}} \left\{ R_{\eta}^2, \hat{s}_{\eta}^{-1} \left( \frac{1}{4\tilde{p}^{[\tilde{j}]} \mu^{[\tilde{j}]2}} (1 - d^{[\tilde{j}]}) \psi_{\Delta}^{[\tilde{j}]2} \right) \right\}$$

and so sufficient conditions are:

$$\left( \beta_{\Delta}^{[j]2} + 2p^{[j]} \tilde{\omega}^{[j]2} \mu_{\Delta}^{[j]2} \right) \leq d^{[j]} u^{[j]2}$$

where:

$$u^{[j]2} = \min_{\tilde{j}:j \in \tilde{\mathbb{P}}^{[\tilde{j}]}} \rho_{(\tilde{\omega}^{[j,\tilde{j}])_{\eta}}^{[j,\tilde{j}]2}} \left\{ R_{\eta}^2, \hat{s}_{\eta}^{-1} \left( \frac{1}{4\tilde{p}^{[\tilde{j}]} \mu^{[\tilde{j}]2}} (1 - d^{[\tilde{j}]}) \psi_{\Delta}^{[\tilde{j}]2} \right) \right\}$$

and the final result follows by setting  $d^{[j]} = \frac{1}{2} \forall j$ .

### C.3 RKHS FORM OF LOCAL MODEL

Our goal in this section is to derive the LiNK kernel from the Local dual model. For reason that will become apparent we find it convenient to work with a slight variant of the local dual feature map with minor re-scaling, namely:

$$\begin{aligned} \tilde{\phi}_{\Delta}^{[\tilde{j},j]}(\mathbf{x}) &= \frac{1}{T_{(\tilde{\omega})_{\eta}}^{[\tilde{j},j]}} \left[ \frac{1}{\eta^k} \left[ \left( \sqrt{\rho_{(\tilde{\omega})_{\eta}}^{[\tilde{j},j]2}} \phi_{\Delta}^{[\tilde{j}]}(\mathbf{x}) \right)^{\otimes l} \right]_{1 \leq l \leq k} \right]_{k \geq 1} \\ \tilde{\mathbf{G}}_{\Delta:i_{\tilde{j}}}^{[\tilde{j},j]}(\mathbf{x}) &= \left[ a_{\left( x_{i_{\tilde{j}}}^{[\tilde{j}]} \right)_k}^{[\tilde{j},j]} \left[ \binom{k}{l} \text{He}_{k-l} \mathbf{G}_{\Delta:i_{\tilde{j}}}^{[\tilde{j}]}(\mathbf{x})^{\otimes l} \right]_{1 \leq l \leq k} \right]_{k \geq 1} \end{aligned} \quad (45)$$

$\forall \tilde{\gamma} \in \tilde{\mathbb{P}}^{[j]}$  and:

$$\begin{aligned} \phi_{\Delta}^{[j]}(\mathbf{x}) &= \sqrt{\frac{1}{\gamma^2+1}} \begin{bmatrix} \frac{1}{\sqrt{2}} \left[ \sqrt{\frac{H^{[j]}}{\tilde{H}^{[j]}}} \frac{1}{\tilde{\omega}^{[j,j]}} \tilde{\mathbf{x}}^{[j,j]} \right]_{\tilde{\gamma} \in \tilde{\mathbb{P}}^{[j]}} \\ \frac{1}{\sqrt{2}} \left[ \sqrt{\frac{H^{[j]}}{\tilde{H}^{[j]}}} \tilde{\phi}_{\Delta}^{[j,j]}(\mathbf{x}) \right]_{\tilde{\gamma} \in \tilde{\mathbb{P}}^{[j]}} \end{bmatrix} \\ \mathbf{G}_{\Delta}^{[j]}(\mathbf{x}) &= \begin{bmatrix} \mathbf{1}_{H^{[j]}}^{\top} \\ \text{diag} \left( \mathbf{I}_{H^{[j]}} \mathbf{1}_{\tilde{H}^{[j]}} \mathbf{1}_{H^{[j]}}^{\top} \right)_{\tilde{\gamma} \in \tilde{\mathbb{P}}^{[j]}} \\ \text{diag} \left( \tilde{\mathbf{G}}_{\Delta}^{[j,j]}(\mathbf{x}) \mathbf{1}_{\tilde{H}^{[j]}} \mathbf{1}_{H^{[j]}}^{\top} \right)_{\tilde{\gamma} \in \tilde{\mathbb{P}}^{[j]}} \end{bmatrix} \end{aligned} \quad (46)$$

We have already demonstrated that the necessary conditions for functions in our space to like in an RKHS, so it remains derive the kernel representing this space. To this end we recall that Reisz representer theory implies that  $\forall \mathbf{x} \in \mathbb{X} \exists$  unique  $\mathbf{K}_{\mathbf{x}} \in \mathcal{F} \times \mathbb{R}^m$  such that  $\langle \mathbf{f}(\mathbf{x}), \mathbf{v} \rangle = \langle \mathbf{f}, \mathbf{K}_{\mathbf{x}} \mathbf{v} \rangle_{\mathcal{H}} \forall \mathbf{v} \in \mathbb{R}^m$ ; from which the kernel may be obtained using:

$$\mathbf{K}(\mathbf{x}, \mathbf{x}') = \left[ \left\langle \mathbf{K}_{\mathbf{x}} \delta_{(i_{D-1})}^{[D-1]}, \mathbf{K}_{\mathbf{x}'} \delta_{(i'_{D-1})}^{[D-1]} \right\rangle_{\mathcal{H}} \right]_{i_{D-1}, i'_{D-1}}$$

where  $\delta_{(k)}^{[j]} = [\delta_{k, i_j}]_{i_j}$ . Thus our first task is to find such a  $\mathbf{K}_{\mathbf{x}}$ . Recall:

$$\Delta \mathbf{f}(\mathbf{x}) = \left( \Psi_{\Delta}^{[D-1]}(\Theta) \odot \mathbf{G}_{\Delta}^{[D-1]}(\mathbf{x}) \right)^{\top} \phi_{\Delta}^{[D-1]}(\mathbf{x})$$

and so:

$$\begin{aligned} \langle \Delta \mathbf{f}(\mathbf{x}), \mathbf{v} \rangle &= \Delta \mathbf{f}(\mathbf{x})^{\top} \mathbf{v} \\ &= \phi_{\Delta}^{[D-1]}(\mathbf{x})^{\top} \left( \Psi_{\Delta}^{[D-1]}(\Theta) \odot \mathbf{G}_{\Delta}^{[D-1]}(\mathbf{x}) \right) \mathbf{v} \\ &= \sum_{k, i'_{D-1}} \phi_{\Delta k}^{[D-1]}(\mathbf{x}) \left( \Psi_{\Delta k, i'_{D-1}}^{[D-1]}(\Theta) G_{\Delta k, i'_{D-1}}^{[D-1]}(\mathbf{x}) \right) v_{i'_{D-1}} \\ &= \sum_{k, i'_{D-1}} \Psi_{\Delta k, i'_{D-1}}^{[D-1]}(\Theta) \left( \phi_{\Delta k}^{[D-1]}(\mathbf{x}) G_{\Delta k, i'_{D-1}}^{[D-1]}(\mathbf{x}) v_{i'_{D-1}} \right) \\ &= \sum_{k, i'_{D-1}} \Psi_{\Delta k, i'_{D-1}}^{[D-1]}(\Theta) \left( \sum_{i''_{D-1}} \delta_{i'_{D-1}, i''_{D-1}} \phi_{\Delta k}^{[D-1]}(\mathbf{x}) G_{\Delta k, i''_{D-1}}^{[D-1]}(\mathbf{x}) v_{i''_{D-1}} \right) \\ &= \langle \mathbf{f}, \mathbf{K}_{\mathbf{x}} \mathbf{v} \rangle \end{aligned}$$

where we have denoted:

$$\begin{aligned} \mathbf{f} &= \left[ \Psi_{\Delta k, i_{D-1}}^{[D-1]}(\Theta) \right]_{(k, i_{D-1})} \\ \mathbf{K}_{\mathbf{x}} &= \left[ \delta_{i_{D-1}, i'_{D-1}} \phi_{\Delta k}^{[D-1]}(\mathbf{x}) G_{\Delta k, i'_{D-1}}^{[D-1]}(\mathbf{x}) \right]_{(k, i_{D-1}), i'_{D-1}} \end{aligned}$$

We may therefore proceed to derive the kernel as follows, letting  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  be the standard inner product:

$$\begin{aligned} K_{i_{D-1}, i'_{D-1}}(\mathbf{x}, \mathbf{x}') &= \left\langle \left\langle \left[ \delta_{\tilde{i}'_{D-1}, \tilde{i}''_{D-1}} \phi_{\Delta k}^{[D-1]}(\mathbf{x}) G_{\Delta k, \tilde{i}''_{D-1}}^{[D-1]}(\mathbf{x}) \right]_{(k, \tilde{i}'_{D-1}), \tilde{i}''_{D-1}} \left[ \delta_{i_{D-1}, \tilde{i}''_{D-1}} \right]_{\tilde{i}''_{D-1}} \right. \right. \\ &\quad \left. \left. \left[ \delta_{\tilde{i}'_{D-1}, \tilde{i}''_{D-1}} \phi_{\Delta k}^{[D-1]}(\mathbf{x}') G_{\Delta k, \tilde{i}''_{D-1}}^{[D-1]}(\mathbf{x}') \right]_{(k, \tilde{i}'_{D-1}), \tilde{i}''_{D-1}} \left[ \delta_{i'_{D-1}, \tilde{i}''_{D-1}} \right]_{\tilde{i}''_{D-1}} \right\rangle_{\mathcal{H}} \right\rangle_{i_{D-1}, i'_{D-1}} \\ &= \left\langle \left\langle \left[ \delta_{\tilde{i}'_{D-1}, i_{D-1}} \phi_{\Delta k}^{[D-1]}(\mathbf{x}) G_{\Delta k, i_{D-1}}^{[D-1]}(\mathbf{x}) \right]_{(k, \tilde{i}'_{D-1})} \right. \right. \\ &\quad \left. \left. \left[ \delta_{\tilde{i}'_{D-1}, i'_{D-1}} \phi_{\Delta k}^{[D-1]}(\mathbf{x}') G_{\Delta k, i'_{D-1}}^{[D-1]}(\mathbf{x}') \right]_{(k, \tilde{i}'_{D-1})} \right\rangle_{\mathcal{H}} \right\rangle_{i_{D-1}, i'_{D-1}} \\ &= \left[ \sum_{k, \tilde{i}'_{D-1}} \delta_{\tilde{i}'_{D-1}, i_{D-1}} \phi_{\Delta k}^{[D-1]}(\mathbf{x}) G_{\Delta k, i_{D-1}}^{[D-1]}(\mathbf{x}) \delta_{\tilde{i}'_{D-1}, i'_{D-1}} \phi_{\Delta k}^{[D-1]}(\mathbf{x}') G_{\Delta k, i'_{D-1}}^{[D-1]}(\mathbf{x}') \right]_{i_{D-1}, i'_{D-1}} \\ &= \left[ \delta_{i_{D-1}, i'_{D-1}} \sum_k \phi_{\Delta k}^{[D-1]}(\mathbf{x}) G_{\Delta k, i_{D-1}}^{[D-1]}(\mathbf{x}) \phi_{\Delta k}^{[D-1]}(\mathbf{x}') G_{\Delta k, i'_{D-1}}^{[D-1]}(\mathbf{x}') \right]_{i_{D-1}, i'_{D-1}} \\ &= \delta_{i_{D-1}, i'_{D-1}} K_{i_{D-1}}^{[D-1]}(\mathbf{x}, \mathbf{x}') \end{aligned}$$

1566

where:

1567

1568

1569

1570

1571

1572

1573

1574

1575

1576

1577

1578

1579

1580

1581

1582

1583

1584

1585

1586

1587

1588

$$\begin{aligned}
K_{i_j}^{[j]}(\mathbf{x}, \mathbf{x}') &= \left\langle \phi_{\Delta}^{[j]}(\mathbf{x}) \odot \mathbf{G}_{\Delta; i_j}^{[j]}(\mathbf{x}), \phi_{\Delta}^{[j]}(\mathbf{x}') \odot \mathbf{G}_{\Delta; i_j}^{[j]}(\mathbf{x}') \right\rangle \\
&= \frac{1}{\gamma^2+1} \left( \gamma^2 + \frac{1}{2} \sum_{\tilde{j} \in \mathbb{P}^{[j]}} \frac{H^{[\tilde{j}]}}{\tilde{H}^{[\tilde{j}]}} \frac{1}{\tilde{\omega}^{[\tilde{j}, j]^2}} \left\langle \tilde{\mathbf{x}}^{[\tilde{j}, j]}, \tilde{\mathbf{x}}'^{[\tilde{j}, j]} \right\rangle + \dots \right. \\
&\quad \left. \dots + \frac{1}{2} \sum_{\tilde{j} \in \mathbb{P}^{[j]}} \frac{H^{[\tilde{j}]}}{\tilde{H}^{[\tilde{j}]}} \sum_{i_{\tilde{j}}} \left\langle \tilde{\phi}_{\Delta}^{[\tilde{j}, j]}(\mathbf{x}) \odot \tilde{\mathbf{G}}_{\Delta; i_{\tilde{j}}}^{[\tilde{j}, j]}(\mathbf{x}), \tilde{\phi}_{\Delta}^{[\tilde{j}, j]}(\mathbf{x}') \odot \tilde{\mathbf{G}}_{\Delta; i_{\tilde{j}}}^{[\tilde{j}, j]}(\mathbf{x}') \right\rangle \right) \\
&= \frac{1}{\gamma^2+1} \left( \gamma^2 + \frac{1}{2} \sum_{\tilde{j} \in \mathbb{P}^{[j]}} \frac{H^{[\tilde{j}]}}{\tilde{H}^{[\tilde{j}]}} \frac{1}{\tilde{\omega}^{[\tilde{j}, j]^2}} \left\langle \tilde{\mathbf{x}}^{[\tilde{j}, j]}, \tilde{\mathbf{x}}'^{[\tilde{j}, j]} \right\rangle + \dots \right. \\
&\quad \left. \dots + \frac{1}{2} \sum_{\tilde{j} \in \mathbb{P}^{[j]}} \frac{H^{[\tilde{j}]}}{\tilde{H}^{[\tilde{j}]}} \frac{1}{T^{[\tilde{j}, j]^2}} \sum_{i_{\tilde{j}}} \sum_{k \geq 1} \frac{a_{i_{\tilde{j}}}^{[\tilde{j}, j]}(x_{i_{\tilde{j}}})^k a_{i_{\tilde{j}}}^{[\tilde{j}, j]}(x'_{i_{\tilde{j}}})^k}{\eta^{2k}} \sum_{1 \leq l \leq k} \binom{k}{l} \text{He}_{k-l}^2 \left( \rho_{(\tilde{\omega})\eta}^{[\tilde{j}, j]^2} \left\langle \phi_{\Delta}^{[\tilde{j}, j]}(\mathbf{x}) \odot \mathbf{G}_{\Delta; i_{\tilde{j}}}^{[\tilde{j}, j]}(\mathbf{x}), \phi_{\Delta}^{[\tilde{j}, j]}(\mathbf{x}') \odot \mathbf{G}_{\Delta; i_{\tilde{j}}}^{[\tilde{j}, j]}(\mathbf{x}') \right\rangle \right)^l \right) \\
&= \frac{1}{\gamma^2+1} \left( \gamma^2 + \sum_{\tilde{j} \in \mathbb{P}^{[j]}} \frac{H^{[\tilde{j}]}}{\tilde{H}^{[\tilde{j}]}} \frac{1}{\tilde{\omega}^{[\tilde{j}, j]^2}} \left\langle \tilde{\mathbf{x}}^{[\tilde{j}, j]}, \tilde{\mathbf{x}}'^{[\tilde{j}, j]} \right\rangle + \dots \right. \\
&\quad \left. \dots + \sum_{\tilde{j} \in \mathbb{P}^{[j]}} \frac{H^{[\tilde{j}]}}{\tilde{H}^{[\tilde{j}]}} \frac{1}{T^{[\tilde{j}, j]^2}} \sum_{i_{\tilde{j}}} \sum_{k \geq 1} \frac{a_{i_{\tilde{j}}}^{[\tilde{j}, j]}(x_{i_{\tilde{j}}})^k a_{i_{\tilde{j}}}^{[\tilde{j}, j]}(x'_{i_{\tilde{j}}})^k}{\eta^{2k}} \sum_{1 \leq l \leq k} \binom{k}{l} \text{He}_{k-l}^2 \left( \rho_{(\tilde{\omega})\eta}^{[\tilde{j}, j]^2} K_{i_{\tilde{j}}}^{[\tilde{j}, j]}(\mathbf{x}, \mathbf{x}') \right)^l \right) \\
&= \frac{1}{\gamma^2+1} \left( \gamma^2 + \sum_{\tilde{j} \in \mathbb{P}^{[j]}} \frac{H^{[\tilde{j}]}}{\tilde{H}^{[\tilde{j}]}} \frac{1}{\tilde{\omega}^{[\tilde{j}, j]^2}} \left\langle \tilde{\mathbf{x}}^{[\tilde{j}, j]}, \tilde{\mathbf{x}}'^{[\tilde{j}, j]} \right\rangle + \sum_{\tilde{j} \in \mathbb{P}^{[j]}} \frac{H^{[\tilde{j}]}}{\tilde{H}^{[\tilde{j}]}} \frac{1}{T^{[\tilde{j}, j]^2}} \sum_{i_{\tilde{j}}} \hat{\tau}_{\eta}^{[\tilde{j}, j]} \left( \rho_{(\tilde{\omega})\eta}^{[\tilde{j}, j]^2} K_{i_{\tilde{j}}}^{[\tilde{j}, j]}(\mathbf{x}, \mathbf{x}'); x_{i_{\tilde{j}}}^{[\tilde{j}, j]}, x'_{i_{\tilde{j}}}^{[\tilde{j}, j]} \right) \right)
\end{aligned}$$

1589

So, noting that this is independent of  $i_j$ , we find that  $\mathbf{f} \in \mathcal{H}_{\mathbf{K}}$ , where  $\mathbf{K} = \mathbf{I}K_{\text{LiNK}}$ , where

1590

 $K_{\text{LiNK}} = K_{\text{LiNK}}^{[D-1]}$  is defined recursively:

1591

1592

1593

1594

with  $K_{\text{LiNK}}^{[-1]}(\mathbf{x}, \mathbf{x}') = 0$ .

1595

1596

Recall that:

1597

1598

1599

$$\hat{\tau}_{\eta}^{[\tilde{j}, j]}(\zeta; \xi, \xi') = \sum_{k=1}^{\infty} \frac{a_{(\xi)k}^{[\tilde{j}, j]} a_{(\xi')k}^{[\tilde{j}, j]}}{\eta^{2k}} \sum_{l=1}^k \binom{k}{l} \text{He}_{k-l}^2 \zeta^l$$

1600

1601

We aim to express:

1602

1603

where  $\hat{\tau}_{\eta, q}^{[\tilde{j}, j]}(\zeta; \xi, \xi') = b(\xi, \xi') \zeta^q$ . As noted previously, if:

1604

1605

1606

then:

1607

1608

1609

$$f^{(q)}(x) = \sum_{k=0}^{\infty} \frac{(k+q)!}{k!} a_{k+q} \sum_{l=0}^{k-q} \binom{k}{l} \text{He}_{k-l} x^l$$

1610

Working toward the general case:

1611

1612

1613

1614

1615

$$\begin{aligned}
\hat{\tau}_{\eta, 1}^{[\tilde{j}, j]}(\zeta; \xi, \xi') &= \zeta \sum_{k=1}^{\infty} \frac{a_{(\xi)k}^{[\tilde{j}, j]} a_{(\xi')k}^{[\tilde{j}, j]}}{\eta^{2k}} \binom{k}{1} \text{He}_{k-1}^2 \\
&= \zeta \sum_{k, k'=0}^{\infty} \delta_{k, k'} \frac{a_{(\xi)k+1}^{[\tilde{j}, j]} a_{(\xi')k'+1}^{[\tilde{j}, j]}}{\eta^{k+1}} \binom{k+1}{1} \text{He}_k \frac{a_{(\xi')k'+1}^{[\tilde{j}, j]} a_{(\xi)k'+1}^{[\tilde{j}, j]}}{\eta^{k'+1}} \binom{k'+1}{1} \text{He}_{k'}
\end{aligned}$$

1616

and:

1617

1618

1619

$$\begin{aligned}
\hat{\tau}_{\eta, 2}^{[\tilde{j}, j]}(\zeta; \xi, \xi') &= \zeta^2 \sum_{k=2}^{\infty} \frac{a_{(\xi)k}^{[\tilde{j}, j]} a_{(\xi')k}^{[\tilde{j}, j]}}{\eta^{2k}} \binom{k}{2} \text{He}_{k-2}^2 \\
&= \zeta^2 \sum_{k, k'=0}^{\infty} \delta_{k, k'} \frac{a_{(\xi)k+2}^{[\tilde{j}, j]} a_{(\xi')k'+2}^{[\tilde{j}, j]}}{\eta^{k+2}} \binom{k+2}{2} \text{He}_k \frac{a_{(\xi')k'+2}^{[\tilde{j}, j]} a_{(\xi)k'+2}^{[\tilde{j}, j]}}{\eta^{k'+2}} \binom{k'+2}{2} \text{He}_{k'}
\end{aligned}$$

and so on, so:

$$\begin{aligned}
\hat{\tau}_{\eta,q}^{[\bar{j},j]}(\zeta; \xi, \xi') &= \zeta^q \sum_{k=q}^{\infty} \frac{a_{(\xi)k}^{[\bar{j},j]} a_{(\xi')k}^{[\bar{j},j]}}{\eta^{qk}} \binom{k}{q}^2 \text{He}_{k-q}^2 \\
&= \zeta^q \sum_{k,k'=0}^{\infty} \delta_{k,k'} \frac{a_{(\xi)k+q}^{[\bar{j},j]}}{\eta^{k+q}} \binom{k+q}{q} \text{He}_k \frac{a_{(\xi')k'+q}^{[\bar{j},j]}}{\eta^{k'+q}} \binom{k'+q}{q} \text{He}_{k'} \\
&= \frac{1}{q!^2} \zeta^q \sum_{k,k'=0}^{\infty} \delta_{k,k'} \frac{a_{(\xi)k+q}^{[\bar{j},j]}}{\eta^{k+q}} \frac{(k+q)!}{k!} \text{He}_k \frac{a_{(\xi')k'+q}^{[\bar{j},j]}}{\eta^{k'+q}} \frac{(k'+q)!}{k'!} \text{He}_{k'} \\
&= \frac{1}{q!^2} \zeta^q \sum_{k,k'=0}^{\infty} \delta_{k,k'} \sum_{l=0}^{k-q} \sum_{l'=0}^{k'-q} \delta_{l,l'} \left( \frac{a_{(\xi)k+q}^{[\bar{j},j]}}{\eta^{k+q}} \frac{(k+q)!}{k!} \binom{k}{l} \text{He}_{k-l} 0^l \right) \left( \frac{a_{(\xi')k'+q}^{[\bar{j},j]}}{\eta^{k'+q}} \frac{(k'+q)!}{k'!} \binom{k'}{l'} \text{He}_{k'-l'} 0^{l'} \right) \\
&= \zeta^q \left( \frac{1}{q!} \sum_{k=0}^{\infty} \frac{a_{(\xi)k+q}^{[\bar{j},j]}}{\eta^{k+q}} \frac{(k+q)!}{k!} \sum_{l=0}^{k-q} \binom{k}{l} \text{He}_{k-l} 0^l \right) \cdot \left( \frac{1}{q!} \sum_{k'=0}^{\infty} \frac{a_{(\xi')k'+q}^{[\bar{j},j]}}{\eta^{k'+q}} \frac{(k'+q)!}{k'!} \sum_{l'=0}^{k'-q} \binom{k'}{l'} \text{He}_{k'-l'} 0^{l'} \right) \\
&= \frac{1}{q!} \bar{\tau}_{\eta}^{[\bar{j},j]}(q)(0; \xi) \frac{1}{q!} \bar{\tau}_{\eta}^{[\bar{j},j]}(q)(0; \xi') \zeta^q
\end{aligned}$$

where  $\cdot$  is the Cauchy product and in the final step we have used Mertens' theorem for Cauchy products, and:

$$\bar{\tau}_{\eta}^{[\bar{j},j]}(\zeta; \xi) = \sum_{k=1}^{\infty} \frac{a_{(\xi)k}^{[\bar{j},j]}}{\eta^k} \sum_{l=1}^k \binom{k}{l} \text{He}_{k-l} \zeta^l$$

Finally, we see that, in the limit  $\eta \rightarrow 1$ :

$$\hat{\tau}_{\eta,q}^{[\bar{j},j]}(\zeta; \xi, \xi') = \frac{1}{q!} \bar{\tau}^{[\bar{j},j]}(q)(0; \xi) \frac{1}{q!} \bar{\tau}^{[\bar{j},j]}(q)(0; \xi') \zeta^q$$

which leads to the simplified form of the LiNK:

$$\begin{aligned}
K_{\text{LiNK}}^{[j]}(\mathbf{x}, \mathbf{x}') &= \frac{1}{\gamma^2+1} \left( \sum_{\bar{j} \in \tilde{\mathbb{P}}^{[j]}} \frac{H^{[\bar{j}]}}{\bar{\omega}^{[\bar{j},j]2}} \Sigma^{[\bar{j},j]}(\mathbf{x}, \mathbf{x}') + \dots \right. \\
&\quad \left. \dots + \sum_{\bar{j} \in \tilde{\mathbb{P}}^{[j]}} \frac{H^{[\bar{j}]}}{H^{[j]}} \frac{1}{T^{[\bar{j},j]2}(\bar{\omega})_{\eta}} \sum_{q=1}^{\infty} \sum_{i_{\bar{j}}} \frac{1}{q!} \bar{\tau}^{[\bar{j},j]}(q)(0; x_{i_{\bar{j}}}^{[\bar{j}]}) \frac{1}{q!} \bar{\tau}^{[\bar{j},j]}(q)(0; x_{i_{\bar{j}}}^{\prime[\bar{j}]}) \left( \rho_{(\bar{\omega})_{\eta}}^{[\bar{j},j]2} K_{\text{LiNK}}^{[j]}(\mathbf{x}, \mathbf{x}') \right)^q \right)
\end{aligned}$$

with  $K_{\text{LiNK}}^{[-1]}(\mathbf{x}, \mathbf{x}') = 0$ . With some cleanup:

$$\begin{aligned}
K_{\text{LiNK}}^{[j]}(\mathbf{x}, \mathbf{x}') &= \frac{1}{\gamma^2+1} \left( \mathbb{E}_{\bar{j} \in \tilde{\mathbb{P}}^{[j]}} \left[ \frac{H^{[\bar{j}]}}{\bar{\omega}^{[\bar{j},j]2}} \Sigma^{[\bar{j},j]}(\mathbf{x}, \mathbf{x}') \right] + \dots \right. \\
&\quad \left. \dots + \mathbb{E}_{\bar{j} \in \tilde{\mathbb{P}}^{[j]}} \left[ \frac{H^{[\bar{j}]}}{T^{[\bar{j},j]2}(\bar{\omega})_{\eta}} \sum_{q=1}^{\infty} \mathbb{E}_{i_{\bar{j}}} \left[ \frac{1}{q!} \bar{\tau}^{[\bar{j},j]}(q)(x_{i_{\bar{j}}}^{[\bar{j}]}) \frac{1}{q!} \bar{\tau}^{[\bar{j},j]}(q)(x_{i_{\bar{j}}}^{\prime[\bar{j}]}) \right] K_{\text{LiNK}}^{[j]}(\mathbf{x}, \mathbf{x}')^q \right] \right)
\end{aligned}$$

with  $K_{\text{LiNK}}^{[-1]}(\mathbf{x}, \mathbf{x}') = 0$ .

## D RADEMACHER COMPLEXITY BOUNDS - PROOF OF THEOREMS 3, 4 AND 8

In this supplementary we prove theorems relating to the Rademacher complexity of neural networks for the global and local models.

**Theorem 3** *The set  $\mathcal{F} = \{f(\cdot; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R} \mid \Theta \in \mathbb{W}\}$  of networks (1) satisfying our assumptions has Rademacher complexity bounded by  $\mathcal{R}_N(\mathcal{F}) \leq \frac{1}{\sqrt{N}} \phi \psi$  (definitions as per Figure 1).*

*Proof.* Let  $\mathcal{F}$  be the set of attainable neural networks (scalar output) and  $\epsilon$  a Rademacher random

variable. Let  $\mathbf{x} \sim \nu$ . Then the Rademacher complexity is bounded as:

$$\begin{aligned}
\mathcal{R}_N(\mathcal{F}) &\triangleq \mathbb{E}_\nu \mathbb{E}_\epsilon \left[ \sup_{f \in \mathcal{F}} \frac{1}{N} \sum_k \epsilon_k f(\mathbf{x}_k) \right] \\
&= \text{Global Dual} \frac{1}{N} \mathbb{E}_\nu \mathbb{E}_\epsilon \left[ \sup_{\Theta \in \mathbb{W}} \sum_k \epsilon_k \langle \Psi(\Theta), \phi(\mathbf{x}_k) \rangle_{\mathbf{g}} \right] \\
&= \frac{1}{N} \mathbb{E}_\nu \mathbb{E}_\epsilon \left[ \sup_{\Theta \in \mathbb{W}} \langle \Psi(\Theta), \sum_k \epsilon_k \phi(\mathbf{x}_k) \rangle_{\mathbf{g}} \right] \\
&\leq \text{Operator norm} \frac{1}{N} \mathbb{E}_\nu \mathbb{E}_\epsilon \left[ \sup_{\Theta \in \mathbb{W}} \|\Psi(\Theta)\|_{\text{He}[\tau]} \sqrt{\|\sum_k \epsilon_k \phi(\mathbf{x}_k)\|_2^2} \right] \\
&\leq \text{Norm Bound} \frac{\psi}{N} \mathbb{E}_\nu \left[ \mathbb{E}_\epsilon \sqrt{\|\sum_k \epsilon_k \phi(\mathbf{x}_k)\|_2^2} \right] \\
&\leq \text{Jensen} \frac{\psi}{N} \mathbb{E}_\nu \left[ \sqrt{\mathbb{E}_\epsilon \|\sum_k \epsilon_k \phi(\mathbf{x}_k)\|_2^2} \right] \\
&= \{\mathbb{E}_\epsilon \epsilon_k \epsilon_l = \delta_{k,l}\} \frac{\psi}{N} \mathbb{E}_\nu \left[ \sqrt{\sum_k \|\phi(\mathbf{x}_k)\|_2^2} \right] \\
&\leq \text{Norm Bound} \frac{\psi}{N} \mathbb{E}_\nu \left[ \sqrt{N\phi^2} \right] = \frac{\phi\psi}{\sqrt{N}}
\end{aligned}$$

□

**Theorem 4** Let  $\mathcal{F} = \{f(\cdot; \Theta) : \mathbb{R}^n \rightarrow \mathbb{R} | \Theta \in \mathbb{W}\}$  be the set of unbiased networks (1) with  $L$ -Lipschitz activations satisfying our assumptions. Then  $\mathcal{R}_N(\mathcal{F}) \leq \max_{\mathcal{S} \in \mathcal{S}} \prod_{j \in \mathcal{S}} L^2 \tilde{p}^{[j]} \mu^{[j]}$ , where  $\mathcal{S}$  is the set of all input-output paths in the network graph.

*Proof.* Observe from Figure 1 that:

$$\tilde{\phi}^{[\tilde{j},j]^2} \tilde{\psi}^{[\tilde{j},j]^2} \leq \frac{\phi^{[j]^2}}{\bar{s}^{[\tilde{j},j]} (\bar{s}^{[\tilde{j},j]-1} (\tilde{\phi}^{[\tilde{j},j]^2})_{\phi^{[j]^2}})} L^2 \tilde{\phi}^{[\tilde{j},j]^2} \tilde{\psi}^{[\tilde{j},j]^2}$$

Note that the numerator of the fractional part is linearly increasing while the denominator is superlinearly increasing, so we may pessimise this bound as:

$$\begin{aligned}
\tilde{\phi}^{[\tilde{j},j]^2} \tilde{\psi}^{[\tilde{j},j]^2} &\leq \frac{\phi^{[j]^2}}{\bar{s}^{[\tilde{j},j]} (\bar{s}^{[\tilde{j},j]-1} (\tilde{\phi}^{[\tilde{j},j]^2})_{\phi^{[j]^2}})} L^2 \tilde{\phi}^{[\tilde{j},j]^2} \tilde{\psi}^{[\tilde{j},j]^2} \\
&\leq \lim_{\phi^{[j]} \rightarrow 0} \frac{\phi^{[j]^2}}{\bar{s}^{[\tilde{j},j]} (\bar{s}^{[\tilde{j},j]-1} (\tilde{\phi}^{[\tilde{j},j]^2})_{\phi^{[j]^2}})} L^2 \tilde{\phi}^{[\tilde{j},j]^2} \tilde{\psi}^{[\tilde{j},j]^2} = \frac{1}{a_0^{[\tilde{j},j]} \bar{s}^{[\tilde{j},j]-1} (\tilde{\phi}^{[\tilde{j},j]^2})} L^2 \tilde{\phi}^{[\tilde{j},j]^2} \tilde{\psi}^{[\tilde{j},j]^2}
\end{aligned}$$

As  $\tilde{\phi}^{[\tilde{j},j]^2}$  is a free parameter here we can let  $\tilde{\phi}^{[\tilde{j},j]^2} \rightarrow 0$ , in which limit  $\tilde{\phi}^{[\tilde{j},j]^2} \tilde{\psi}^{[\tilde{j},j]^2} \leq L^2 \tilde{\psi}^{[\tilde{j},j]^2} \leq L^2 p^{[j]} \mu^{[j]^2} \sum_{\tilde{j} \in \tilde{\mathbb{P}}^{[j]}} \tilde{\phi}^{[\tilde{j},j]^2} \tilde{\psi}^{[\tilde{j},j]^2}$ . The result follows recursively, then upper bounding with the path-wise product. □

**Theorem 8** The set  $\mathcal{F}_\Delta = \{\Delta f(\cdot; \Delta\Theta) : \mathbb{R}^n \rightarrow \mathbb{R} | \Delta\Theta \in \mathbb{W}_\Delta\}$  of change in neural-network operation satisfying (20) has Rademacher complexity  $\mathcal{R}_N(\mathcal{F}) \leq \frac{1}{\sqrt{N}} \phi_\Delta \psi_\Delta$  (defined in Figure 2).

*Proof.* Let  $\mathcal{F}_\Delta$  be the set of attainable changes in neural networks (scalar output) and  $\epsilon$  a Rademacher

random variable. Let  $\mathbf{x} \sim \nu$ . Then the Rademacher complexity is bounded as:

$$\begin{aligned}
\mathcal{R}_N(\mathcal{F}_\Delta) &\triangleq \mathbb{E}_\nu \mathbb{E}_\epsilon \left[ \sup_{\Delta f \in \mathcal{F}} \frac{1}{N} \sum_k \epsilon_k \Delta f(\mathbf{x}_k) \right] \\
&\stackrel{\text{Local Dual}}{=} \frac{1}{N} \mathbb{E}_\nu \mathbb{E}_\epsilon \left[ \sup_{\Delta \Theta \in \mathbb{W}_\Delta} \sum_k \epsilon_k \langle \Psi_\Delta(\Delta \Theta), \phi_\Delta(\mathbf{x}_k) \rangle_{\mathbf{g}} \right] \\
&= \frac{1}{N} \mathbb{E}_\nu \mathbb{E}_\epsilon \left[ \sup_{\Delta \Theta \in \mathbb{W}_\Delta} \langle \Psi_\Delta(\Delta \Theta), \sum_k \epsilon_k \phi_\Delta(\mathbf{x}_k) \rangle_{\mathbf{g}} \right] \\
&\leq \text{Cauchy Schwarz} \frac{1}{N} \mathbb{E}_\nu \mathbb{E}_\epsilon \left[ \sup_{\Delta \Theta \in \mathbb{W}_\Delta} \|\Psi_\Delta(\Delta \Theta)\|_2 \sqrt{\|\sum_k \epsilon_k \phi_\Delta(\mathbf{x}_k)\|_2^2} \right] \\
&\leq \text{Norm Bound} \frac{\psi_\Delta}{N} \mathbb{E}_\nu \left[ \mathbb{E}_\epsilon \sqrt{\|\sum_k \epsilon_k \phi_\Delta(\mathbf{x}_k)\|_2^2} \right] \\
&\leq \text{Jensen} \frac{\psi_\Delta}{N} \mathbb{E}_\nu \left[ \sqrt{\mathbb{E}_\epsilon \|\sum_k \epsilon_k \phi_\Delta(\mathbf{x}_k)\|_2^2} \right] \\
&= \{\mathbb{E}_\epsilon \epsilon_k \epsilon_l = \delta_{k,l}\} \frac{\psi_\Delta}{N} \mathbb{E}_\nu \left[ \sqrt{\sum_k \|\phi_\Delta(\mathbf{x}_k)\|_2^2} \right] \\
&\leq \text{Norm Bound} \frac{\psi_\Delta}{N} \mathbb{E}_\nu \left[ \sqrt{N \phi_\Delta^2} \right] = \frac{\phi_\Delta \psi_\Delta}{\sqrt{N}}
\end{aligned}$$

□

1728  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1780  
1781