# Optimal Approximation Rate of ReLU Networks in terms of Width and Depth[*]

Zuowei Shen[†]    Haizhao Yang[‡]    Shijun Zhang[§]

**Abstract**

This paper concentrates on the approximation power of deep feed-forward neural networks in terms of width and depth. It is proved by construction that ReLU networks with width $\mathcal{O}\big(\max\{d\lfloor N^{1/d}\rfloor, N+2\}\big)$ and depth $\mathcal{O}(L)$ can approximate a Hölder continuous function on $[0,1]^d$ with an approximation rate $\mathcal{O}\big(\lambda\sqrt{d}(N^2L^2\ln N)^{-\alpha/d}\big)$, where $\alpha\in(0,1]$ and $\lambda>0$ are Hölder order and constant, respectively. Such a rate is optimal up to a constant in terms of width and depth separately, while existing results are only nearly optimal without the logarithmic factor in the approximation rate. More generally, for an arbitrary continuous function $f$ on $[0,1]^d$, the approximation rate becomes $\mathcal{O}\big(\sqrt{d}\,\omega_f\big((N^2L^2\ln N)^{-1/d}\big)\big)$, where $\omega_f(\cdot)$ is the modulus of continuity. We also extend our analysis to any continuous function $f$ on a bounded set. Particularly, if ReLU networks with depth 31 and width $\mathcal{O}(N)$ are used to approximate one-dimensional Lipschitz continuous functions on $[0,1]$ with a Lipschitz constant $\lambda>0$, the approximation rate in terms of the total number of parameters, $W=\mathcal{O}(N^2)$, becomes $\mathcal{O}(\frac{\lambda}{W\ln W})$, which has not been discovered in the literature for fixed-depth ReLU networks.

**Key words**. Deep ReLU Networks; Optimal Approximation; VC-dimension; Bit Extraction; Hölder Continuity.

## 1 Introduction

Over the past few decades, the expressiveness of neural networks has been widely studied from many points of view, e.g., in terms of combinatorics [18], topology [4], Vapnik-Chervonenkis (VC) dimension [3,8,21], fat-shattering dimension [1,12], information theory [20], classical approximation theory [2,5,7,9,13,15,22,22,23,24,25,26,28,29], optimization [10,11,19]. The error analysis of neural networks consists of three parts: the approximation error, the optimization error, and the generalization error. This paper focuses on the approximation error for ReLU networks.

The approximation errors of feed-forward neural networks with various activation functions have been studied for different types of functions, e.g., smooth functions [6, 14, 15, 16, 27], piecewise smooth functions [20], band-limited functions [17], continuous functions [23, 24, 25, 28]. In [23], it was shown that a ReLU network with width $C_1(d) \cdot N$ and depth $C_2(d) \cdot L$ can attain an approximation error $C_3(d) \cdot \omega_f(N^{-2/d}L^{-2/d})$ to approximate a continuous function $f$ on $[0,1]^d$, where $C_1(d)$, $C_2(d)$, and $C_3(d)$ are three constants in $d$ with explicit formulas to specify their values, and $\omega_f(\cdot)$ is the modulus of continuity of $f \in C([0,1]^d)$ defined via

$$\omega_f(r) := \sup\left\{|f(\boldsymbol{x}) - f(\boldsymbol{y})| : \boldsymbol{x}, \boldsymbol{y} \in [0,1]^d, \ \|\boldsymbol{x} - \boldsymbol{y}\|_2 \le r\right\}, \quad \text{for any } r \ge 0.$$

Such an approximation rate is optimal in terms of $N$ and $L$ up to a logarithmic term and the corresponding optimal approximation theory is still open. To address this open problem, we provide a constructive proof in this paper to show that ReLU networks of width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ can approximate an arbitrary continuous function $f$ on $[0,1]^d$ with an optimal approximation error $\mathcal{O}\left(\sqrt{d}\,\omega_f\left((N^2L^2\ln N)^{-\alpha/d}\right)\right)$ in terms of $N$ and $L$. As shown by our main result, Theorem 1.1 below, the approximation rate obtained here admits explicit formulas to specify its prefactors when $\omega_f(\cdot)$ is known.

**Theorem 1.1.** *Given a continuous function $f \in C([0,1]^d)$, for any $N \in \mathbb{N}^+$, $L \in \mathbb{N}^+$, and $p \in [1, \infty]$, there exists a function $\phi$ implemented by a ReLU network with width $C_1 \max\left\{d\lfloor N^{1/d}\rfloor, N+2\right\}$ and depth $11L + C_2$ such that*

$$\|f - \phi\|_{L^p([0,1]^d)} \le 131\sqrt{d}\,\omega_f\left(\left(N^2L^2\log_3(N+2)\right)^{-1/d}\right),$$

*where $C_1 = 16$ and $C_2 = 18$ if $p \in [1, \infty)$; $C_1 = 3^{d+3}$ and $C_2 = 18 + 2d$ if $p = \infty$.*

Note that $3^{d+3}\max\left\{d\lfloor N^{1/d}\rfloor, N+2\right\} \le 3^{d+3}\max\left\{dN, 3N\right\} \le 3^{d+4}dN$. Given any $\widetilde{N}, \widetilde{L} \in \mathbb{N}^+$ with $\widetilde{N} \ge 3^{d+4}d$ and $\widetilde{L} \ge 29 + 2d$, there exist $N, L \in \mathbb{N}^+$ such that

$$3^{d+4}dN \le \widetilde{N} < 3^{d+4}d(N+1) \quad \text{and} \quad 11L + 18 + 2d \le \widetilde{L} < 11(L+1) + 18 + 2d.$$

If follows that

$$N \ge \frac{N+1}{3} > \frac{\widetilde{N}}{3^{d+5}d} \quad \text{and} \quad L \ge \frac{L+1}{2} > \frac{1}{2} \cdot \frac{\widetilde{L} - 18 - 2d}{11} = \frac{\widetilde{L} - 18 - 2d}{22}.$$

Then we have an immediate corollary of Theorem 1.1.

**Corollary 1.2.** *Given a continuous function $f \in C([0,1]^d)$, for any $\widetilde{N}, \widetilde{L} \in \mathbb{N}^+$ with $\widetilde{N} \ge 3^{d+4}d$ and $\widetilde{L} \ge 29 + 2d$, there exists a function $\phi$ implemented by a ReLU network with width $\widetilde{N}$ and depth $\widetilde{L}$ such that*

$$\|f - \phi\|_{L^\infty([0,1]^d)} \le 131\sqrt{d}\,\omega_f\left(\left(\left(\tfrac{\widetilde{N}}{3^{d+5}d}\right)^2\left(\tfrac{\widetilde{L}-18-2d}{22}\right)^2\log_3\left(\tfrac{\widetilde{N}}{3^{d+5}d}+2\right)\right)^{-1/d}\right).$$

As a special case of Theorem 1.1 for explicit error characterization, let us take Hölder continuous functions as an example. Let $\text{Hölder}([0,1]^d, \alpha, \lambda)$ denote the space of Hölder continuous functions on $[0,1]^d$ of order $\alpha \in (0,1]$ with a Hölder constant $\lambda > 0$. We have an immediate corollary of Theorem 1.1 as follows.

**Corollary 1.3.** *Given a Hölder continuous function $f \in \text{Hölder}([0,1]^d, \alpha, \lambda)$, for any $N \in \mathbb{N}^+$, $L \in \mathbb{N}^+$, and $p \in [1, \infty]$, there exists a function $\phi$ implemented by a ReLU network with width $C_1 \max\left\{ d \lfloor N^{1/d} \rfloor, N + 2 \right\}$ and depth $11L + C_2$ such that*

$$\|f - \phi\|_{L^p([0,1]^d)} \leq 131\lambda\sqrt{d}\big(N^2 L^2 \log_3(N+2)\big)^{-\alpha/d},$$

*where $C_1 = 16$ and $C_2 = 18$ if $p \in [1, \infty)$; $C_1 = 3^{d+3}$ and $C_2 = 18 + 2d$ if $p = \infty$.*

To better illustrate the importance of our theory, we summarize our key contributions as follows.

(1) Upper bound: We provide a quantitative and non-asymptotic approximation rate $131\sqrt{d}\,\omega_f\big(\big(N^2 L^2 \log_3(N+2)\big)^{-1/d}\big)$ in terms of width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ for any $f \in C([0,1]^d)$ in Theorem 1.1.

   (1.1) This approximation error analysis can be extended to $f \in C(E)$ for any $E \subseteq [-R, R]^d$ with $R > 0$ as we shall see later in Theorem 2.5.

   (1.2) In the case of one-dimensional Lipschitz continuous functions on $[0,1]$ with a constant $\lambda > 0$, the approximation rate in Theorem 1.1 becomes $\mathcal{O}(\frac{\lambda}{W \ln W})$ for ReLU networks with 31 hidden layers and $\mathcal{O}(W)$ parameters via setting $L = 1$ and $W = \mathcal{O}(N^2)$ therein. To the best of our knowledge, the approximation rate $\mathcal{O}(\frac{\lambda}{W \ln W})$ is better than existing known results using fixed-depth ReLU networks to approximate Lipschitz continuous functions on $[0,1]$.

(2) Lower bound: Through the VC-dimension bounds of ReLU networks given in [8], we show, in Section 2.3, that the approximation rate $131\lambda\sqrt{d}\big(N^2 L^2 \log_3(N+2)\big)^{-\alpha/d}$ in terms of width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ for $\text{Hölder}([0,1]^d, \alpha, \lambda)$ is optimal as follows.

   (2.1) When the width is fixed, both the approximation upper and lower bounds take the form of $CL^{-2\alpha/d}$ for a positive constant $C$.

   (2.2) When the depth is fixed, both the approximation upper and lower bounds take the form of $C(N^2 \ln N)^{-\alpha/d}$ for a positive constant $C$.

We would like to point out that if $N$ and $L$ vary simultaneously, the rate is optimal in the $N$-$L$ plane except for a small region as shown in Figure 1. See Section 2.3 for a detailed discussion. The earlier result in [23] provides a nearly optimal approximation error that has a gap (a logarithmic term) between the lower and upper bounds. It is technically challenging to match the upper bound with the lower bound. Compared to the nearly optimal rate $19\lambda\sqrt{d}N^{-2\alpha/d}L^{-2\alpha/d}$ for Hölder continuous functions in $\text{Hölder}([0,1]^d, \alpha, \lambda)$ in [23], this paper achieves the optimal rate $131\lambda\sqrt{d}\big(N^2 L^2 \log_3(N+2)\big)^{-\alpha/d}$ using more technical and sophisticated construction. For example, a novel bit extraction technique different to that in [3] is proposed, and new ReLU networks are constructed to approximate step functions more efficiently than those in [23]. The optimal result obtained in this paper could also be extended to other functions spaces, leading to better understanding of deep network approximation.
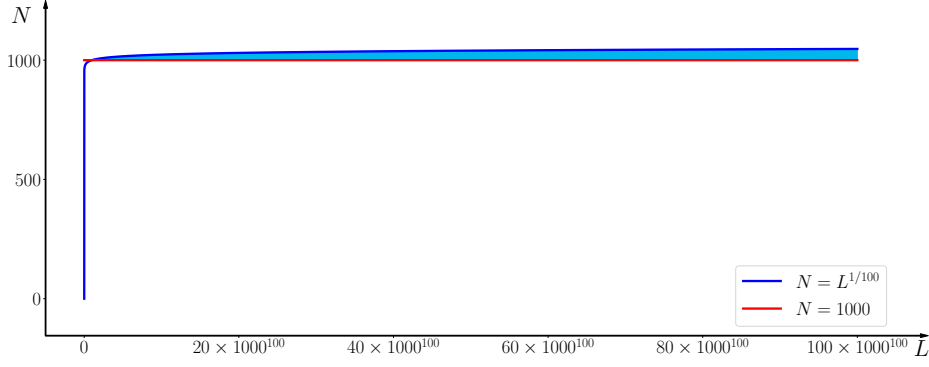
3

Figure 1: Our rate is optimal in terms of width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ simultaneously except for the region marked in cyan characterized by $\{(N, L) \in \mathbb{N}^2 : C_1 \le N \le L^{C_2}\}$, where $C_i = C_i(\alpha, d)$ for $i = 1, 2$ are two positive constants. This figure is an example for $C_1 = 1000$ and $C_2 = 1/100$.

The error analysis of deep learning is to estimate approximation, generalization, and optimization errors. Here, we give a brief discussion, the interested reader can find more details in [15]. Let $\phi(\boldsymbol{x}; \boldsymbol{\theta})$ denote a function computed by a network parameterized with $\boldsymbol{\theta}$. Given a target function $f$, the final goal is to find the expected risk minimizer

$$\boldsymbol{\theta}_{\mathcal{D}} \coloneqq \arg\min_{\boldsymbol{\theta}} R_{\mathcal{D}}(\boldsymbol{\theta}), \quad \text{where } R_{\mathcal{D}}(\boldsymbol{\theta}) \coloneqq \mathbb{E}_{\boldsymbol{x} \sim U(\mathcal{X})} \left[\ell(\phi(\boldsymbol{x}; \boldsymbol{\theta}), f(\boldsymbol{x}))\right]$$

with a loss function $\ell(\cdot, \cdot)$ and an unknown data distribution $U(\mathcal{X})$.

In practice, for given samples $\{(\boldsymbol{x}_i, f(\boldsymbol{x}_i))\}_{i=1}^n$, the goal of supervised learning is to identify the empirical risk minimizer

$$\boldsymbol{\theta}_{\mathcal{S}} \coloneqq \arg\min_{\boldsymbol{\theta}} R_{\mathcal{S}}(\boldsymbol{\theta}), \quad \text{where } R_{\mathcal{S}}(\boldsymbol{\theta}) \coloneqq \frac{1}{n} \sum_{i=1}^n \ell\big(\phi(\boldsymbol{x}_i; \boldsymbol{\theta}), f(\boldsymbol{x}_i)\big).$$

In fact, one could only get a numerical minimizer $\boldsymbol{\theta}_{\mathcal{N}}$ via a numerical optimization method. The discrepancy between the target function $f$ and the learned function $\phi(\boldsymbol{x}; \boldsymbol{\theta}_{\mathcal{N}})$ is measured by $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}})$, which is bounded by

$$R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}}) \le \underbrace{R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})}_{\text{Approximation error}} + \underbrace{[R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}})]}_{\text{Optimization error}} + \underbrace{[R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}})] + [R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{D}}) - R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})]}_{\text{Generalization error}}.$$

This paper deals with the approximation error of ReLU networks for continues functions and gives an upper bound of $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})$ which is optimal up to a constant. Note that the approximation error analysis given here is independent of data samples and deep learning algorithms. However, the analysis of optimization and generalization errors do depend on data samples, deep learning algorithms, models, etc.

The rest of this paper is organized as follows. In Section 2, we prove Theorem 1.1 by assuming Theorem 2.1 is true, show the optimality of Theorem 1.1, and extend our analysis to continuous functions defined on any bounded set. Next, Theorem 2.1 is proved in Section 3 based on Proposition 3.1 and 3.2, the proofs of which can be found in Section 4. Finally, Section 5 concludes this paper with a short discussion.

4

# 2 Theoretical analysis

In this section, we first prove Theorem 1.1 and discuss its optimality. Next, we extend our analysis to general continuous functions defined on any bounded set in $\mathbb{R}^d$. Notations throughout this paper is summarized in Section 2.1.

## 2.1 Notations

Let us summarize all basic notations used in this paper as follows.

- Matrices are denoted by bold uppercase letters. For instance, $\boldsymbol{A} \in \mathbb{R}^{m \times n}$ is a real matrix of size $m \times n$, and $\boldsymbol{A}^T$ denotes the transpose of $\boldsymbol{A}$. Vectors are denoted as bold lowercase letters. For example, $\boldsymbol{v} = [v_1, \cdots, v_d]^T = \begin{bmatrix} v_1 \\ \vdots \\ v_d \end{bmatrix} \in \mathbb{R}^d$ is a column vector with $\boldsymbol{v}(i) = v_i$ being the $i$-th element. Besides, "[" and "]" are used to partition matrices (vectors) into blocks, e.g., $\boldsymbol{A} = \begin{bmatrix} \boldsymbol{A}_{11} & \boldsymbol{A}_{12} \\ \boldsymbol{A}_{21} & \boldsymbol{A}_{22} \end{bmatrix}$.

- For any $p \in [1, \infty)$, the $p$-norm (or $\ell^p$-norm) of a vector $\boldsymbol{x} = [x_1, x_2, \cdots, x_d]^T \in \mathbb{R}^d$ is defined by
$$\|\boldsymbol{x}\|_p \coloneqq \left( |x_1|^p + |x_2|^p + \cdots + |x_d|^p \right)^{1/p}.$$

- For any $x \in \mathbb{R}$, let $\lfloor x \rfloor \coloneqq \max\{n : n \le x, \ n \in \mathbb{Z}\}$ and $\lceil x \rceil \coloneqq \min\{n : n \ge x, \ n \in \mathbb{Z}\}$.

- Assume $\boldsymbol{n} \in \mathbb{N}^d$, then $f(\boldsymbol{n}) = \mathcal{O}(g(\boldsymbol{n}))$ means that there exists positive $C$ independent of $\boldsymbol{n}$, $f$, and $g$ such that $f(\boldsymbol{n}) \le Cg(\boldsymbol{n})$ when all entries of $\boldsymbol{n}$ go to $+\infty$.

- For any $\theta \in [0, 1)$, suppose its binary representation is $\theta = \sum_{\ell=1}^{\infty} \theta_\ell 2^{-\ell}$ with $\theta_\ell \in \{0, 1\}$, we introduce a special notation $\mathrm{bin}\, 0.\theta_1\theta_2\cdots\theta_L$ to denote the $L$-term binary representation of $\theta$, i.e., $\mathrm{bin}\, 0.\theta_1\theta_2\cdots\theta_L \coloneqq \sum_{\ell=1}^{L} \theta_\ell 2^{-\ell}$.

- Let $\mu(\cdot)$ denote the Lebesgue measure.

- Let $1_S$ be the characteristic function on a set $S$, i.e., $1_S$ is equal to 1 on $S$ and 0 outside $S$.

- Let $|S|$ denote the size of a set $S$, i.e., the number of all elements in $S$.

- The set difference of two sets $A$ and $B$ is denoted by $A \backslash B \coloneqq \{x : x \in A, \ x \notin B\}$.

- Given any $K \in \mathbb{N}^+$ and $\delta \in (0, \frac{1}{K})$, define a trifling region $\Omega([0, 1]^d, K, \delta)$ of $[0, 1]^d$ as
$$\Omega([0, 1]^d, K, \delta) \coloneqq \bigcup_{j=1}^{d} \left\{ \boldsymbol{x} = [x_1, x_2, \cdots, x_d]^T \in [0, 1]^d : x_j \in \bigcup_{k=1}^{K-1} \left( \frac{k}{K} - \delta, \frac{k}{K} \right) \right\}. \qquad (2.1)$$

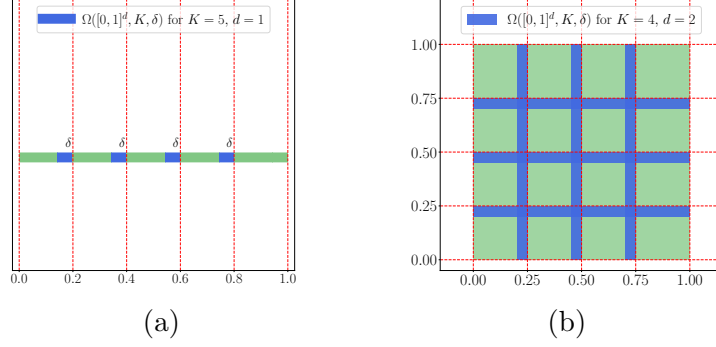In particular, $\Omega([0, 1]^d, K, \delta) = \varnothing$ if $K = 1$. See Figure 2 for two examples of trifling regions.

5

Figure 2: Two examples of trifling regions. (a) $K = 5, d = 1$. (b) $K = 4, d = 2$.

- Let Hölder$([0,1]^d, \alpha, \lambda)$ denote the space of Hölder continuous functions on $[0,1]^d$ of order $\alpha \in (0,1]$ with a Hölder constant $\lambda > 0$.

- For a continuous piecewise linear function $f(x)$, the $x$ values where the slope changes are typically called **breakpoints**.

- Let $\mathrm{CPwL}(\mathbb{R}, n)$ denote the space that consists of all continuous piecewise linear functions with at most $n$ breakpoints on $\mathbb{R}$.

- Let $\sigma : \mathbb{R} \to \mathbb{R}$ denote the rectified linear unit (ReLU), i.e. $\sigma(x) = \max\{0, x\}$. With a slight abuse of notation, we define $\sigma : \mathbb{R}^d \to \mathbb{R}^d$ as $\sigma(\boldsymbol{x}) = \begin{bmatrix} \max\{0, x_1\} \\ \vdots \\ \max\{0, x_d\} \end{bmatrix}$ for any $\boldsymbol{x} = [x_1, \cdots, x_d]^T \in \mathbb{R}^d$.

- We will use $\mathcal{NN}$ to denote a function implemented by a ReLU network for short and use Python-type notations to specify a class of functions implemented by ReLU networks with several conditions, e.g., $\mathcal{NN}(c_1; c_2; \cdots; c_m)$ is a set of functions implemented by ReLU networks satisfying $m$ conditions given by $\{c_i\}_{1 \le i \le m}$, each of which may specify the number of inputs (#input), the number of outputs (#output), the total number of neurons in all hidden layers (#neuron), the number of hidden layers (depth), the total number of parameters (#parameter), and the width in each hidden layer (widthvec), the maximum width of all hidden layers (width), etc. For example, if $\phi \in \mathcal{NN}(\text{#input} = 2; \text{widthvec} = [100, 100]; \text{#output} = 1)$, then $\phi$ is a functions satisfies

  - $\phi$ maps from $\mathbb{R}^2$ to $\mathbb{R}$.
  - $\phi$ can be implemented by a ReLU network with two hidden layers and the number of nodes in each hidden layer is 100.

- For any function $\phi \in \mathcal{NN}(\text{#input} = d; \text{widthvec} = [N_1, N_2, \cdots, N_L]; \text{#output} = 1)$, if we set $N_0 = d$ and $N_{L+1} = 1$, then the architecture of the network implementing $\phi$ can be briefly described as follows:

$$\boldsymbol{x} = \widetilde{\boldsymbol{h}}_0 \xrightarrow[\mathcal{L}_0]{\boldsymbol{W}_0, \, \boldsymbol{b}_0} \boldsymbol{h}_1 \xrightarrow{\sigma} \widetilde{\boldsymbol{h}}_1 \cdots \xrightarrow[\mathcal{L}_{L-1}]{\boldsymbol{W}_{L-1}, \, \boldsymbol{b}_{L-1}} \boldsymbol{h}_L \xrightarrow{\sigma} \widetilde{\boldsymbol{h}}_L \xrightarrow[\mathcal{L}_L]{\boldsymbol{W}_L, \, \boldsymbol{b}_L} \boldsymbol{h}_{L+1} = \phi(\boldsymbol{x}),$$

6

where $\boldsymbol{W}_i \in \mathbb{R}^{N_{i+1} \times N_i}$ and $\boldsymbol{b}_i \in \mathbb{R}^{N_{i+1}}$ are the weight matrix and the bias vector in the $i$-th affine linear transform $\mathcal{L}_i$ in $\phi$, respectively, i.e.,

$$\boldsymbol{h}_{i+1} = \boldsymbol{W}_i \cdot \widetilde{\boldsymbol{h}}_i + \boldsymbol{b}_i =: \mathcal{L}_i(\widetilde{\boldsymbol{h}}_i), \quad \text{for } i = 0, 1, \cdots, L,$$

and

$$\widetilde{\boldsymbol{h}}_i = \sigma(\boldsymbol{h}_i), \quad \text{for } i = 1, \ldots, L.$$

In particular, $\phi$ can be represented in a form of function compositions as follows.

$$\phi = \mathcal{L}_L \circ \sigma \circ \mathcal{L}_{L-1} \circ \sigma \circ \cdots \circ \sigma \circ \mathcal{L}_1 \circ \sigma \circ \mathcal{L}_0,$$

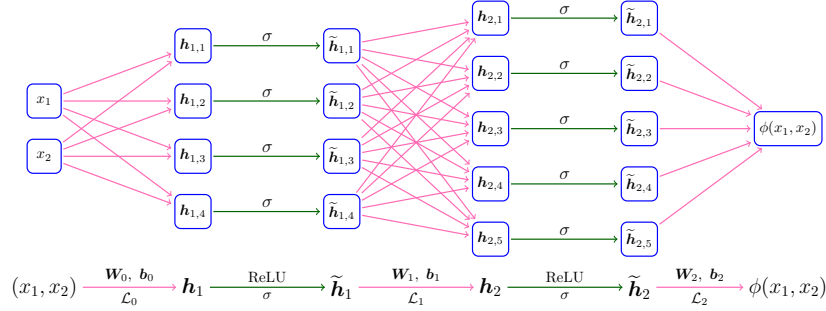which has been illustrated in Figure 3.



Figure 3: An example of a ReLU network with width 5 and depth 2.

- The expression "a network with width $N$ and depth $L$" means

    – The maximum width of this network for all **hidden** layers is no more than $N$.

    – The number of **hidden** layers of this network is no more than $L$.

## 2.2 Proof of Theorem 1.1

The key point is to construct piecewise constant functions to approximate continuous functions in the proof. However, it is impossible to construct a piecewise constant function implemented by a ReLU network due to the continuity of ReLU networks. Thus, we introduce the trifling region $\Omega([0,1]^d, K, \delta)$, defined in Equation (2.1), and use ReLU networks to implement piecewise constant functions outside the trifling region. To prove Theorem 1.1, we first introduce a weaker variant of Theorem 1.1, showing how to construct ReLU networks to pointwisely approximate continuous functions except for the trifling region.

**Theorem 2.1.** *Given a function $f \in C([0,1]^d)$, for any $N \in \mathbb{N}^+$ and $L \in \mathbb{N}^+$, there exists a function $\phi$ implemented by a ReLU network with width $\max\{8d\lfloor N^{1/d}\rfloor + 3d, 16N + 30\}$ and depth $11L + 18$ such that $\|\phi\|_{L^\infty(\mathbb{R}^d)} \leq |f(\boldsymbol{0})| + \omega_f(\sqrt{d})$ and*

$$|f(\boldsymbol{x}) - \phi(\boldsymbol{x})| \leq 130\sqrt{d}\,\omega_f\big((N^2 L^2 \log_3(N+2))^{-1/d}\big), \quad \text{for any } \boldsymbol{x} \in [0,1]^d \backslash \Omega([0,1]^d, K, \delta),$$

*where $K = \lfloor N^{1/d}\rfloor^2 \lfloor L^{1/d}\rfloor^2 \lfloor\lfloor\log_3(N+2)\rfloor^{1/d}\rfloor$ and $\delta$ is an arbitrary number in $(0, \frac{1}{3K}]$.*

7

With Theorem 2.1 that will be proved in Section 3, we can easily prove Theorem 1.1 for the case $p \in [1, \infty)$. To attain the rate in $L^\infty$-norm, we need to control the approximation error in the trifling region. To this end, we introduce a theorem to deal with the approximation inside the trifling region $\Omega([0,1]^d, K, \delta)$.

**Theorem 2.2** (Theorem 3.7 of [29] or Theorem 2.1 of [15])*. Given any $\varepsilon > 0$, $N, L, K \in \mathbb{N}^+$, and $\delta \in (0, \frac{1}{3K}]$, assume $f$ is a continuous function in $C([0,1]^d)$ and $\widetilde{\phi}$ can be implemented by a ReLU network with width $N$ and depth $L$. If*

$$|f(\boldsymbol{x}) - \widetilde{\phi}(\boldsymbol{x})| \le \varepsilon, \quad \text{for any } \boldsymbol{x} \in [0,1]^d \backslash \Omega([0,1]^d, K, \delta),$$

*then there exists a function $\phi$ implemented by a new ReLU network with width $3^d(N+4)$ and depth $L + 2d$ such that*

$$|f(\boldsymbol{x}) - \phi(\boldsymbol{x})| \le \varepsilon + d \cdot \omega_f(\delta), \quad \text{for any } \boldsymbol{x} \in [0,1]^d.$$

Now we are ready to prove Theorem 1.1 by assuming Theorem 2.1 is true, which will be proved later in Section 3.

*Proof of Theorem 1.1.* We may assume $f$ is not a constant function since it is a trivial case. Then $\omega_f(r) > 0$ for any $r > 0$. Let us first consider the case $p \in [1, \infty)$. Set $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor^2 \lfloor \lfloor \log_3(N+2) \rfloor^{1/d} \rfloor$ and choose a small $\delta \in (0, \frac{1}{3K}]$ such that

$$
\begin{aligned}
K d \delta \big(2|f(\boldsymbol{0})| + 2\omega_f(\sqrt{d})\big)^p &= \lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor^2 \lfloor \lfloor \log_3(N+2) \rfloor^{1/d} \rfloor d \delta \big(2|f(\boldsymbol{0})| + 2\omega_f(\sqrt{d})\big)^p \\
&\le \left( \omega_f\big( (N^2 L^2 \log_3(N+2))^{-1/d} \big) \right)^p.
\end{aligned}
$$

By Theorem 2.1, there exists a function $\phi$ implemented by a ReLU network with width

$$\max\big\{ 8d\lfloor N^{1/d} \rfloor + 3d, 16N + 30 \big\} \le 16 \max\big\{ d\lfloor N^{1/d} \rfloor, N + 2 \big\}$$

and depth $11L + 18$ such that $\|\phi\|_{L^\infty(\mathbb{R}^d)} \le |f(\boldsymbol{0})| + \omega_f(\sqrt{d})$ and

$$|f(\boldsymbol{x}) - \phi(\boldsymbol{x})| \le 130\sqrt{d}\,\omega_f\big( (N^2 L^2 \log_3(N+2))^{-1/d} \big), \quad \text{for any } \boldsymbol{x} \in [0,1]^d \backslash \Omega([0,1]^d, K, \delta),$$

It follows from $\mu(\Omega([0,1]^d, K, \delta)) \le Kd\delta$ and $\|f\|_{L^\infty([0,1]^d)} \le |f(\boldsymbol{0})| + \omega_f(\sqrt{d})$ that

$$
\begin{aligned}
\|f - \phi\|_{L^p([0,1]^d)}^p &= \int_{\Omega([0,1]^d, K, \delta)} |f(\boldsymbol{x}) - \phi(\boldsymbol{x})|^p \mathrm{d}\boldsymbol{x} + \int_{[0,1]^d \backslash \Omega([0,1]^d, K, \delta)} |f(\boldsymbol{x}) - \phi(\boldsymbol{x})|^p \mathrm{d}\boldsymbol{x} \\
&\le Kd\delta \big(2|f(\boldsymbol{0})| + 2\omega_f(\sqrt{d})\big)^p + \left( 130\sqrt{d}\,\omega_f\big( (N^2 L^2 \log_3(N+2))^{-1/d} \big) \right)^p \\
&\le \left( \omega_f\big( (N^2 L^2 \log_3(N+2))^{-1/d} \big) \right)^p + \left( 130\sqrt{d}\,\omega_f\big( (N^2 L^2 \log_3(N+2))^{-1/d} \big) \right)^p \\
&\le \left( 131\sqrt{d}\,\omega_f\big( (N^2 L^2 \log_3(N+2))^{-1/d} \big) \right)^p.
\end{aligned}
$$

Hence, $\|f - \phi\|_{L^p([0,1]^d)} \le 131\sqrt{d}\,\omega_f\left(\left(N^2 L^2 \log_3(N+2)\right)^{-1/d}\right)$.

Next, let us discuss the case $p = \infty$. Set $K = \lfloor N^{1/d}\rfloor^2 \lfloor L^{1/d}\rfloor^2 \lfloor \lfloor\log_3(N+2)\rfloor^{1/d}\rfloor$ and choose a small $\delta \in (0, \frac{1}{3K}]$ such that

$$d \cdot \omega_f(\delta) \le \omega_f\left(\left(N^2 L^2 \log_3(N+2)\right)^{-1/d}\right).$$

By Theorem 2.1, there exists a function $\widetilde{\phi}$ implemented by a ReLU network with width $\max\left\{8d\lfloor N^{1/d}\rfloor + 3d,\, 16N + 30\right\}$ and depth $11L + 18$ such that

$$|f(\boldsymbol{x}) - \widetilde{\phi}(\boldsymbol{x})| \le 130\sqrt{d}\,\omega_f\left(\left(N^2 L^2 \log_3(N+2)\right)^{-1/d}\right) =: \varepsilon,$$

for any $\boldsymbol{x} \in [0,1]^d\backslash\Omega([0,1]^d, K, \delta)$. By Theorem 2.2, there exists a function $\phi$ implemented by a ReLU network with width

$$3^d\left(\max\left\{8d\lfloor N^{1/d}\rfloor + 3d,\, 16N + 30\right\} + 4\right) \le 3^{d+3}\max\left\{d\lfloor N^{1/d}\rfloor,\, N + 2\right\}$$

and depth $11L + 18 + 2d$ such that

$$|f(\boldsymbol{x}) - \phi(\boldsymbol{x})| \le \varepsilon + d \cdot \omega_f(\delta) \le 131\sqrt{d}\,\omega_f\left(\left(N^2 L^2 \log_3(N+2)\right)^{-1/d}\right), \quad \text{for any } \boldsymbol{x} \in [0,1]^d.$$

So we finish the proof. $\qquad\square$

## 2.3 Optimality

This section will show that the approximation rates in Theorem 1.1 and Corollary 1.3 are optimal and there is no room to improve for the function class Hölder$([0,1]^d, \alpha, \lambda)$. Therefore, the approximation rate for the whole continuous functions space in terms of width and depth in Theorem 1.1 cannot be improved. A typical method to characterize the optimal approximation theory of neural networks is to study the connection between the approximation error and Vapnik–Chervonenkis (VC) dimension [15, 23, 27, 28, 29]. This method relies on the VC-dimension upper bound given in [8]. In this paper, we adopt this method with several modifications to simplify the proof.

Let us first present the definitions of VC-dimension and related concepts. Let $H$ be a class of functions mapping from a general domain $\mathcal{X}$ to $\{0,1\}$. We say $H$ shatters the set $\{\boldsymbol{x}_1, \boldsymbol{x}_2, \cdots, \boldsymbol{x}_m\} \subseteq \mathcal{X}$ if

$$\left|\left\{\left[h(\boldsymbol{x}_1), h(\boldsymbol{x}_2), \cdots, h(\boldsymbol{x}_m)\right]^T \in \{0,1\}^m : h \in H\right\}\right| = 2^m,$$

where $|\cdot|$ denotes the size of a set. This equation means, given any $\theta_i \in \{0,1\}$ for $i = 1, 2, \cdots, m$, there exists $h \in H$ such that $h(\boldsymbol{x}_i) = \theta_i$ for all $i$. For general a function set $\mathscr{F}$ mapping from $\mathcal{X}$ to $\mathbb{R}$, we say $\mathscr{F}$ shatters $\{\boldsymbol{x}_1, \boldsymbol{x}_2, \cdots, \boldsymbol{x}_m\} \subseteq \mathcal{X}$ if $\mathcal{T} \circ \mathscr{F}$ does, where

$$\mathcal{T}(t) := \begin{cases} 1, & t \ge 0, \\ 0, & t < 0 \end{cases} \quad \text{and} \quad \mathcal{T} \circ \mathscr{F} := \{\mathcal{T} \circ f : f \in \mathscr{F}\}.$$

For any $m \in \mathbb{N}^+$, we define the growth function of $H$ as

$$\Pi_H(m) := \max_{\boldsymbol{x}_1, \boldsymbol{x}_2, \cdots, \boldsymbol{x}_m \in \mathcal{X}} \left|\left\{\left[h(\boldsymbol{x}_1), h(\boldsymbol{x}_2), \cdots, h(\boldsymbol{x}_m)\right]^T \in \{0,1\}^m : h \in H\right\}\right|.$$

**Definition 2.3** (VC-dimension). Let $H$ be a class of functions from $\mathcal{X}$ to $\{0,1\}$. The VC-dimension of $H$, denoted by $\mathrm{VCDim}(H)$, is the size of the largest shattered set, namely,

$$\mathrm{VCDim}(H) \coloneqq \sup\{m \in \mathbb{N}^+ : \Pi_H(m) = 2^m\}$$

if $\{m \in \mathbb{N}^+ : \Pi_H(m) = 2^m\}$ is not empty. In the case of $\{m \in \mathbb{N}^+ : \Pi_H(m) = 2^m\} = \varnothing$, we may define $\mathrm{VCDim}(H) = 0$.

Let $\mathscr{F}$ be a class of functions from $\mathcal{X}$ to $\mathbb{R}$. The VC-dimension of $\mathscr{F}$, denoted by $\mathrm{VCDim}(\mathscr{F})$, is defined by $\mathrm{VCDim}(\mathscr{F}) \coloneqq \mathrm{VCDim}(\mathcal{T} \circ \mathscr{F})$, where

$$\mathcal{T}(t) \coloneqq \begin{cases} 1, & t \geq 0, \\ 0, & t < 0 \end{cases} \quad \text{and} \quad \mathcal{T} \circ \mathscr{F} \coloneqq \{\mathcal{T} \circ f : f \in \mathscr{F}\}.$$

In particular, the expression "VC-dimension of a network (architecture)" means the VC-dimension of the function set that consists of all functions implemented by this network (architecture).

We remark that one may also define $\mathrm{VCDim}(\mathscr{F})$ as $\mathrm{VCDim}(\mathscr{F}) \coloneqq \mathrm{VCDim}(\widetilde{\mathcal{T}} \circ \mathscr{F})$, where

$$\widetilde{\mathcal{T}}(t) \coloneqq \begin{cases} 1, & t > 0, \\ 0, & t \leq 0 \end{cases} \quad \text{and} \quad \widetilde{\mathcal{T}} \circ \mathscr{F} \coloneqq \{\widetilde{\mathcal{T}} \circ f : f \in \mathscr{F}\}.$$

Note that function spaces generated by networks are closed under linear transformation. Thus, these two definitions of VC-dimension are equivalent.

The theorem below, similar to Theorem 4.17 of [29], reveals the connection between VC-dimension and approximation rate.

**Theorem 2.4.** *Assume $\mathscr{F}$ is a set of functions mapping from $[0,1]^d$ to $\mathbb{R}$. For any $\varepsilon > 0$, if $\mathrm{VCDim}(\mathscr{F}) \geq 1$ and*

$$\inf_{\phi \in \mathscr{F}} \|\phi - f\|_{L^\infty([0,1]^d)} \leq \varepsilon, \quad \text{for any } f \in \mathrm{Hölder}([0,1]^d, \alpha, 1), \tag{2.2}$$

*then $\mathrm{VCDim}(\mathscr{F}) \geq (9\varepsilon)^{-d/\alpha}$.*

This theorem demonstrates the connection between VC-dimension of $\mathscr{F}$ and the approximation rate using elements of $\mathscr{F}$ to approximate functions in $\mathrm{Hölder}([0,1]^d, \alpha, \lambda)$. To be precise, the VC-dimension of $\mathscr{F}$ determines an approximation rate lower bound $\mathrm{VCDim}(\mathscr{F})^{-\alpha/d}/9$, which is the best possible approximation rate. Denote the best approximation error of functions in $\mathrm{Hölder}([0,1]^d, \alpha, 1)$ approximated by ReLU networks with width $N$ and depth $L$ as

$$\mathcal{E}_{\alpha,d}(N, L) \coloneqq \sup_{f \in \mathrm{Hölder}([0,1]^d, \alpha, 1)} \left( \inf_{\phi \in \mathcal{NN}(\mathrm{width} \leq N; \mathrm{depth} \leq L)} \|\phi - f\|_{L^\infty([0,1]^d)} \right),$$

We have three remarks listed below.

(i) A large VC-dimension cannot guarantee a good approximation rate. For example, it is easy to verify that

$$\mathrm{VCDim}\Big(\{f : f(x) = \cos(ax), \ a \in \mathbb{R}\}\Big) = \infty.$$

10

However, functions in $\{f : f(x) = \cos(ax), \ a \in \mathbb{R}\}$ cannot approximate Hölder continuous functions well.

(ii) A large VC-dimension is necessary for a good approximation rate, because the best possible approximation rate is controlled by an expression of VC-dimension, as shown in Theorem 2.4. For example, Theorem 6 and 8 of [8] implies that

$$\text{VCDim}\big(\mathcal{NN}(\text{width} \le N; \ \text{depth} \le L)\big) \le \min\Big\{\mathcal{O}\big(N^2 L^2 \ln(NL)\big), \mathcal{O}(N^3 L^2)\Big\},$$

deducing

$$\underbrace{C_1(\alpha, d)\Big(\min\{N^2 L^2 \ln(NL), N^3 L^2\}\Big)^{-\alpha/d}}_{\text{implied by Theorem 2.4}} \le \mathcal{E}_{\alpha, d}(N, L) \le \underbrace{C_2(\alpha, d)\big(N^2 L^2 \ln N\big)^{-\alpha/d}}_{\text{implied by Corollary 1.2 and 1.3}}, \quad (2.3)$$

where $C_1(\alpha, d)$ and $C_2(\alpha, d)$ are two positive constants determined by $s, d$, and $C_2(s, d)$ can be explicitly expressed.

- When $L = L_0$ is fixed, Equation (2.3) implies

$$C_1(\alpha, d, L_0)(N^2 \ln N)^{-\alpha/d} \le \mathcal{E}_{\alpha, d}(N, L_0) \le C_2(\alpha, d, L_0)(N^2 \ln N)^{-\alpha/d},$$

where $C_1(\alpha, d, L_0)$ and $C_2(\alpha, d, L_0)$ are two positive constants determined by $\alpha, d, L_0$.

- When $N = N_0$ is fixed, Equation (2.3) implies

$$C_1(\alpha, d, N_0)L^{-2\alpha/d} \le \mathcal{E}_{\alpha, d}(N_0, L) \le C_2(\alpha, d, N_0)L^{-2\alpha/d},$$

where $C_1(\alpha, d, N_0)$ and $C_2(\alpha, d, N_0)$ are two positive constants determined by $\alpha, d, N_0$.

- It is easy to verify that Equation (2.3) is tight except for the following region

$$\big\{(N, L) \in \mathbb{N}^2 : C_3(\alpha, d) \le N \le L^{C_4(\alpha, d)}\big\},$$

$C_3 = C_3(\alpha, d)$ and $C_4 = C_4(\alpha, d)$ are two positive constants. See Figure 1 for an illustration for the case $C_3 = 1000$ and $C_4 = 1/100$.

Finally, let us present the detailed proof of Theorem 2.4.

*Proof of Theorem 2.4.* Recall that the VC-dimension of a function set is defined as the size of the largest set of points that this class of functions can shatter. So our goal is to find a subset of $\mathscr{F}$ to shatter $\mathcal{O}(\varepsilon^{-d/\alpha})$ points in $[0, 1]^d$, which can be divided into two steps.

- Construct $\{f_\chi : \chi \in \mathscr{B}\} \subseteq \text{Hölder}([0, 1]^d, \alpha, 1)$ that scatters $\mathcal{O}(\varepsilon^{-d/\alpha})$ points, where $\mathscr{B}$ is a set defined later.

- Design $\phi_\chi \in \mathscr{F}$, for each $\chi \in \mathscr{B}$, based on $f_\chi$ and Equation (2.2) such that $\{\phi_\chi : \chi \in \mathscr{F}\} \subseteq \mathscr{F}$ also shatters $\mathcal{O}(\varepsilon^{-d/\alpha})$ points.

11

327 The details of these two steps can be found below.

**Step** 1: Construct $\{f_\chi : \chi \in \mathscr{B}\} \subseteq \text{Hölder}([0,1]^d, \alpha, 1)$ that scatters $\mathcal{O}(\varepsilon^{-d/\alpha})$ points.

329 We may assume $\varepsilon \le 2/9$ since the case $\varepsilon > 2/9$ is trivial. In fact, $\varepsilon > 2/9$ implies

$$\text{VCDim}(\mathscr{F}) \ge 1 \ge 1/2 \ge 2^{-d/\alpha} > (9\varepsilon)^{-d/\alpha}.$$

331 Let $K = \lfloor (9\varepsilon/2)^{-1/\alpha} \rfloor \in \mathbb{N}^+$ and divide $[0,1]^d$ into $K^d$ non-overlapping sub-cubes $\{Q_{\boldsymbol{\beta}}\}_{\boldsymbol{\beta}}$
332 as follows:

$$Q_{\boldsymbol{\beta}} \coloneqq \Big\{ \boldsymbol{x} = [x_1, x_2, \cdots, x_d]^T \in [0,1]^d : x_i \in \big[\tfrac{\beta_i}{K}, \tfrac{\beta_i+1}{K}\big], \ i = 1, 2, \cdots, d \Big\},$$

334 for any index vector $\boldsymbol{\beta} = [\beta_1, \beta_2, \cdots, \beta_d]^T \in \{0, 1, \cdots, K-1\}^d$.

335 Define a function $\zeta_Q$ on $[0,1]^d$ corresponding to $Q = Q(\boldsymbol{x}_0, \eta) \subseteq [0,1]^d$ such that:

336 - $\zeta_Q(\boldsymbol{x}_0) = (\eta/2)^\alpha/2$;

337 - $\zeta_Q(\boldsymbol{x}) = 0$ for any $\boldsymbol{x} \notin Q\backslash\partial Q$, where $\partial Q$ is the boundary of $Q$;

338 - $\zeta_Q$ is linear on the line that connects $\boldsymbol{x}_0$ and $\boldsymbol{x}$ for any $\boldsymbol{x} \in \partial Q$.

339 Define
340
$$\mathscr{B} \coloneqq \big\{ \chi : \chi \text{ is a map from } \{0, 1, \cdots, K-1\}^d \text{ to } \{-1, 1\} \big\}.$$

341 For each $\chi \in \mathscr{B}$, we define

$$f_\chi(\boldsymbol{x}) \coloneqq \sum_{\boldsymbol{\beta} \in \{0,1,\cdots,K-1\}^d} \chi(\boldsymbol{\beta}) \zeta_{Q_{\boldsymbol{\beta}}}(\boldsymbol{x}),$$

343 where $\zeta_{Q_{\boldsymbol{\beta}}}(\boldsymbol{x})$ is the associated function introduced just above. It is easy to check that
344 $\{f_\chi : \chi \in \mathscr{B}\} \subseteq \text{Hölder}([0,1]^d, \alpha, 1)$ can shatter $K^d = \mathcal{O}(\varepsilon^{-d/\alpha})$ points in $[0,1]^d$.

345 **Step** 2: Construct $\{\phi_\chi : \chi \in \mathscr{B}\}$ that also scatters $\mathcal{O}(\varepsilon^{-d/\alpha})$ points.

346 By Equation (2.2), for each $\chi \in \mathscr{B}$, there exists $\phi_\chi \in \mathscr{F}$ such that

$$\|\phi_\chi - f_\chi\|_{L^\infty([0,1]^d)} \le \varepsilon + \varepsilon/81.$$

348 Let $\mu(\cdot)$ denote the Lebesgue measure of a set. Then, for each $\chi \in \mathscr{B}$, there exists
349 $\mathcal{H}_\chi \subseteq [0,1]^d$ with $\mu(\mathcal{H}_\chi) = 0$ such that

$$|\phi_\chi(\boldsymbol{x}) - f_\chi(\boldsymbol{x})| \le \tfrac{82}{81}\varepsilon, \quad \text{for any } \boldsymbol{x} \in [0,1]\backslash\mathcal{H}_\chi.$$

351 Set $\mathcal{H} = \cup_{\chi \in \mathscr{B}} \mathcal{H}_\chi$, then we have $\mu(\mathcal{H}) = 0$ and

$$|\phi_\chi(\boldsymbol{x}) - f_\chi(\boldsymbol{x})| \le \tfrac{82}{81}\varepsilon, \quad \text{for any } \chi \in \mathscr{B} \text{ and } \boldsymbol{x} \in [0,1]\backslash\mathcal{H}. \tag{2.4}$$

353 Since $Q_{\boldsymbol{\beta}}$ has a sidelength $\frac{1}{K} = \frac{1}{\lfloor (9\varepsilon/2)^{-1/\alpha} \rfloor}$, we have, for each $\boldsymbol{\beta} \in \{0, 1, \cdots, K-1\}^d$ and
354 any $\boldsymbol{x} \in \frac{1}{10}Q_{\boldsymbol{\beta}}$[①],

$$|f_\chi(\boldsymbol{x})| = |\zeta_{Q_{\boldsymbol{\beta}}}(\boldsymbol{x})| \ge \tfrac{9}{10}|\zeta_{Q_{\boldsymbol{\beta}}}(\boldsymbol{x}_{Q_{\boldsymbol{\beta}}})| = \tfrac{9}{10}\big(\tfrac{1}{2\lfloor (9\varepsilon/2)^{-1/\alpha}\rfloor}\big)^\alpha/2 \ge \tfrac{81}{80}\varepsilon, \tag{2.5}$$

---
[①] $\frac{1}{10}Q_{\boldsymbol{\beta}}$ denotes the closed cube whose sidelength is 1/10 of that of $Q_{\boldsymbol{\beta}}$ and which shares the same center of $Q_{\boldsymbol{\beta}}$.

where $\boldsymbol{x}_{Q_{\boldsymbol{\beta}}}$ is the center of $Q_{\boldsymbol{\beta}}$.

Note that $(\frac{1}{10}Q_{\boldsymbol{\beta}})\backslash\mathcal{H}$ is not empty, since $\mu\big((\frac{1}{10}Q_{\boldsymbol{\beta}})\backslash\mathcal{H}\big) > 0$ for each $\boldsymbol{\beta} \in \{0, 1, \cdots, K - 1\}^d$. Together with Equation (2.4) and (2.5), there exists $\boldsymbol{x}_{\boldsymbol{\beta}} \in (\frac{1}{10}Q_{\boldsymbol{\beta}})\backslash\mathcal{H}$ such that, for each $\boldsymbol{\beta} \in \{0, 1, \cdots, K - 1\}^d$ and each $\chi \in \mathscr{B}$,

$$|f_\chi(\boldsymbol{x}_{\boldsymbol{\beta}})| \geq \tfrac{81}{80}\varepsilon > \tfrac{82}{81}\varepsilon \geq |f_\chi(\boldsymbol{x}_{\boldsymbol{\beta}}) - \phi_\chi(\boldsymbol{x}_{\boldsymbol{\beta}})|,$$

Hence, $f_\chi(\boldsymbol{x}_{\boldsymbol{\beta}})$ and $\phi_\chi(\boldsymbol{x}_{\boldsymbol{\beta}})$ have the same sign for each $\chi \in \mathscr{B}$ and $\boldsymbol{\beta} \in \{0, 1, \cdots, K - 1\}^d$. Then $\{\phi_\chi : \chi \in \mathscr{B}\}$ shatters $\big\{\boldsymbol{x}_{\boldsymbol{\beta}} : \boldsymbol{\beta} \in \{0, 1, \cdots, K - 1\}^d\big\}$ since $\{f_\chi : \chi \in \mathscr{B}\}$ shatters $\big\{\boldsymbol{x}_{\boldsymbol{\beta}} : \boldsymbol{\beta} \in \{0, 1, \cdots, K - 1\}^d\big\}$. Therefore,

$$\mathrm{VCDim}(\mathscr{F}) \geq \mathrm{VCDim}\big(\{\phi_\chi : \chi \in \mathscr{B}\}\big) \geq K^d = \lfloor (9\varepsilon/2)^{-1/\alpha} \rfloor^d \geq (9\varepsilon)^{-d/\alpha}, \qquad (2.6)$$

where the last inequality comes from the fact $\lfloor x \rfloor \geq x/2 \geq x/(2^{1/\alpha})$ for any $x \in [1, \infty)$ and $\alpha \in (0, 1]$. So we finish the proof. $\qquad\square$

## 2.4 Approximation in irregular domain

We extend our analysis to general continuous functions defined on any irregular bounded set in $\mathbb{R}^d$. The key idea is to extend the target function to a hypercube while preserving the modulus of continuity. For a general set $E \subseteq \mathbb{R}^d$, the modulus of continuity of $f \in C(E)$ is defined via

$$\omega_f^E(r) := \sup\big\{|f(\boldsymbol{x}) - f(\boldsymbol{y})| : \boldsymbol{x}, \boldsymbol{y} \in E, \ \|\boldsymbol{x} - \boldsymbol{y}\|_2 \leq r\big\}, \quad \text{for any } r \geq 0.$$

In particular, $\omega_f(\cdot)$ is short of $\omega_f^E(\cdot)$ in the case of $E = [0, 1]^d$. Then, Theorem 1.1 can be generalized to $f \in C(E)$ for any bounded set $E \subseteq [-R, R]^d$ with $R > 0$, as shown in the following theorem.

**Theorem 2.5.** *Given a continuous function $f \in C(E)$ with $E \subseteq [-R, R]^d$ and $R > 0$, for any $N \in \mathbb{N}^+$, $L \in \mathbb{N}^+$, and $p \in [1, \infty]$, there exists a function $\phi$ implemented by a ReLU network with width $C_1 \max\big\{d\lfloor N^{1/d}\rfloor, N + 2\big\}$ and depth $11L + C_2$ such that*

$$\|f - \phi\|_{L^p(E)} \leq 131(2R)^{d/p}\sqrt{d}\,\omega_f^E\Big(2R\big(N^2 L^2 \log_3(N + 2)\big)^{-1/d}\Big),$$

*where $C_1 = 16$ and $C_2 = 18$ if $p \in [1, \infty)$; $C_1 = 3^{d+3}$ and $C_2 = 18 + 2d$ if $p = \infty$.*

*Proof.* Given any $f \in C(E)$, by Lemma 4.2 of [23] via setting $S = [-R, R]^d$, there exists $g \in C([-R, R]^d)$ such that

- $g(\boldsymbol{x}) = f(\boldsymbol{x})$ for any $\boldsymbol{x} \in E \subseteq S = [-R, R]^d$;

- $\omega_g^S(r) = \omega_f^E(r)$ for any $r \geq 0$.

Define
$$\widetilde{g}(\boldsymbol{x}) := g(2R\boldsymbol{x} - R), \quad \text{for any } \boldsymbol{x} \in [0, 1]^d.$$

13

By applying Theorem 1.1 to $\widetilde{g} \in C([0,1]^d)$, there exists a function $\widetilde{\phi}$ implemented by a ReLU network with width $C_1 \max\left\{ d\lfloor N^{1/d}\rfloor,\, N+2 \right\}$ and depth $11L + C_2$ such that

$$\|\widetilde{\phi} - \widetilde{g}\|_{L^p([0,1]^d)} \le 131\sqrt{d}\,\omega_{\widetilde{g}}\Big(\big(N^2 L^2 \log_3(N+2)\big)^{-1/d}\Big),$$

where $C_1 = 16$ and $C_2 = 18$ if $p \in [1,\infty)$; $C_1 = 3^{d+3}$ and $C_2 = 18 + 2d$ if $p = \infty$.

Recall that $f(\boldsymbol{x}) = g(\boldsymbol{x}) = \widetilde{g}\big(\frac{\boldsymbol{x}+R}{2R}\big)$ for any $\boldsymbol{x} \in E \subseteq S = [-R,R]^d$ and

$$\omega_{\widetilde{g}}(r) = \omega_g^S(2Rr) = \omega_f^E(2Rr), \quad \text{for any } r \ge 0.$$

Define $\phi(\boldsymbol{x}) \coloneqq \widetilde{\phi}\big(\frac{\boldsymbol{x}+R}{2R}\big) = \widetilde{\phi} \circ \mathcal{L}(\boldsymbol{x})$ for any $\boldsymbol{x} \in \mathbb{R}^d$, where $\mathcal{L}$ is an affine linear map given by $\mathcal{L}(\boldsymbol{x}) = \frac{\boldsymbol{x}+R}{2R}$. Clearly, $\phi$ can be implemented by a ReLU network with width $C_1 \max\left\{ d\lfloor N^{1/d}\rfloor,\, N+2 \right\}$ and depth $11L + C_2$, where $C_1 = 16$ and $C_2 = 18$ if $p \in [1,\infty)$; $C_1 = 3^{d+3}$ and $C_2 = 18 + 2d$ if $p = \infty$. Moreover, for any $\boldsymbol{x} \in E \subseteq S = [-R,R]^d$, we have $\frac{\boldsymbol{x}+R}{2R} \in [0,1]^d$, implying

$$
\begin{aligned}
\|\phi - f\|_{L^p(E)} = \|\phi - g\|_{L^p(E)} &= \|\widetilde{\phi}\circ\mathcal{L} - \widetilde{g}\circ\mathcal{L}\|_{L^p(E)} \\
&\le \|\widetilde{\phi}\circ\mathcal{L} - \widetilde{g}\circ\mathcal{L}\|_{L^p([-R,R]^d)} = (2R)^{d/p}\|\widetilde{\phi} - \widetilde{g}\|_{L^p([0,1]^d)} \\
&\le 131(2R)^{d/p}\sqrt{d}\,\omega_{\widetilde{g}}\Big(\big(N^2 L^2 \log_3(N+2)\big)^{-1/d}\Big) \\
&= 131(2R)^{d/p}\sqrt{d}\,\omega_f^E\Big(2R\big(N^2 L^2 \log_3(N+2)\big)^{-1/d}\Big).
\end{aligned}
$$

With the discussion above, we have proved Theorem 2.5. $\qquad\square$

# 3    Proof of Theorem 2.1

We will prove Theorem 2.1 in this section. We first present the key ideas in Section 3.1. The detailed proof is presented in Section 3.2, based on two propositions in Section 3.1, the proofs of which can be found in Section 4.

## 3.1    Key ideas of proving Theorem 2.1

Given an arbitrary $f \in C([0,1]^d)$, our goal is to construct an almost piecewise constant function $\phi$ implemented by a ReLU network to approximate $f$ well. To this end, we introduce a piecewise constant function $f_p \approx f$ serving as an intermediate approximant in our construction in the sense that

$$f \approx f_p \text{ on } [0,1]^d \quad \text{and} \quad f_p \approx \phi \text{ on } [0,1]^d \backslash \Omega([0,1]^d, K, \delta).$$

The approximation in $f \approx f_p$ is a simple and standard technique in constructive approximation. The most technical part is to design a ReLU network with the desired width and depth to implement a function $\phi$ with $\phi \approx f_p$ outside $\Omega([0,1]^d, K, \delta)$. See Figure 4 for an illustration. The introduction of the trifling region is to ease the construction of $\phi$, which is a continuous piecewise linear function, to approximate the discontinuous function $f_p$ by removing the difficulty near discontinuous points, essentially smoothing $f_p$ by restricting the approximation domain in $[0,1]^d \backslash \Omega([0,1]^d, K, \delta)$.

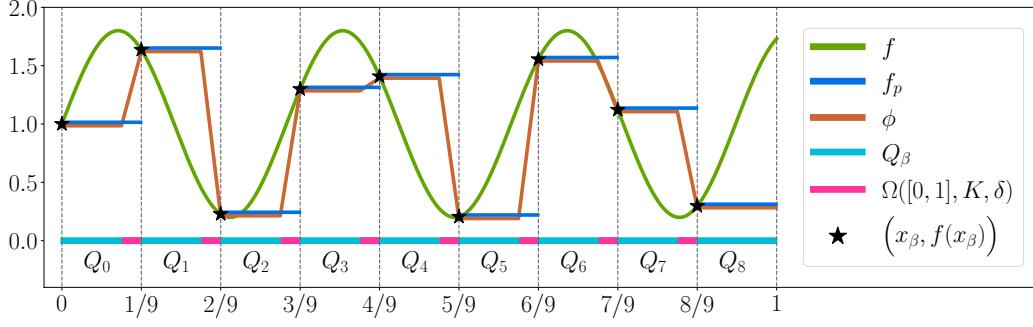Now let us discuss the detailed steps of construction.

14

Figure 4: An illustration of $f$, $f_p$, $\phi$, $x_\beta$, $Q_\beta$, and the trifling region $\Omega([0,1]^d, K, \delta)$ in the one-dimensional case for $\beta \in \{0, 1, \cdots, K-1\}^d$, where $K = N^2 L^2 \log_3(N+2)$ and $d = 1$ with $N = 1$ and $L = 3$. $f$ is the target function; $f_p$ is the piecewise constant function approximating $f$; $\phi$ is a function, implemented by a ReLU network, approximating $f$; and $x_\beta$ is a representative of $Q_\beta$. The measure of $\Omega([0,1]^d, K, \delta)$ can be arbitrarily small as we shall see in the proof of Theorem 1.1.

(1) First, divide $[0,1]^d$ into a union of important regions $\{Q_\beta\}_\beta$ and the trifling region $\Omega([0,1]^d, K, \delta)$, where each $Q_\beta$ is associated with a representative $\boldsymbol{x_\beta} \in Q_\beta$ such that $f(\boldsymbol{x_\beta}) = f_p(\boldsymbol{x_\beta})$ for each index vector $\boldsymbol{\beta} \in \{0, 1, \ldots, K-1\}^d$, where $K = \mathcal{O}((N^2 L^2 \ln N)^{1/d})$ is the partition number per dimension (see Figure 6 for examples for $d = 1$ and $d = 2$).

(2) Next, we design a vector function $\boldsymbol{\Phi}_1(\boldsymbol{x})$ constructed via

$$\boldsymbol{\Phi}_1(\boldsymbol{x}) = \big[\phi_1(x_1), \phi_1(x_2), \cdots, \phi_1(x_d)\big]^T$$

to project the whole cube $Q_\beta$ to a $d$-dimensional index $\boldsymbol{\beta}$ for each $\boldsymbol{\beta}$, where each one-dimensional function $\phi_1$ is a step function implemented by a ReLU network.

(3) The third step is to solve a point fitting problem. To be precise, we construct a function $\phi_2$ implemented by a ReLU network to map $\boldsymbol{\beta}$ approximately to $f_p(\boldsymbol{x_\beta}) = f(\boldsymbol{x_\beta})$. Then $\phi_2 \circ \boldsymbol{\Phi}_1(\boldsymbol{x}) = \phi_2(\boldsymbol{\beta}) \approx f_p(\boldsymbol{x_\beta}) = f(\boldsymbol{x_\beta})$ for any $\boldsymbol{x} \in Q_\beta$ and each $\boldsymbol{\beta}$, implying $\phi := \phi_2 \circ \boldsymbol{\Phi}_1 \approx f_p \approx f$ on $[0,1]^d \backslash \Omega([0,1]^d, K, \delta)$. We would like to point out that we only need to care about the values of $\phi_2$ at a set of points $\{0, 1, \cdots, K-1\}^d$ in the construction of $\phi_2$ according to our design $\phi = \phi_2 \circ \boldsymbol{\Phi}_1$ as illustrated in Figure 5. Therefore, it is not necessary to care about the values of $\phi_2$ sampled outside the set $\{0, 1, \cdots, K-1\}^d$, which is a key point to ease the design of a ReLU network to implement $\phi_2$ as we shall see later.

Finally, we discuss how to implement $\boldsymbol{\Phi}_1$ and $\phi_2$ by deep ReLU networks with width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ using two propositions as we shall prove in Section 4.2 and 4.3 later. We first show how to construct a ReLU network with the desired width and depth by Proposition 3.1 to implement a one-dimensional step function $\phi_1$. Then $\boldsymbol{\Phi}_1$ can be attained via defining

$$\boldsymbol{\Phi}_1(\boldsymbol{x}) = \big[\phi_1(x_1), \phi_1(x_2), \cdots, \phi_1(x_d)\big]^T, \quad \text{for any } \boldsymbol{x} = [x_1, x_2, \cdots, x_d]^T \in \mathbb{R}^d.$$
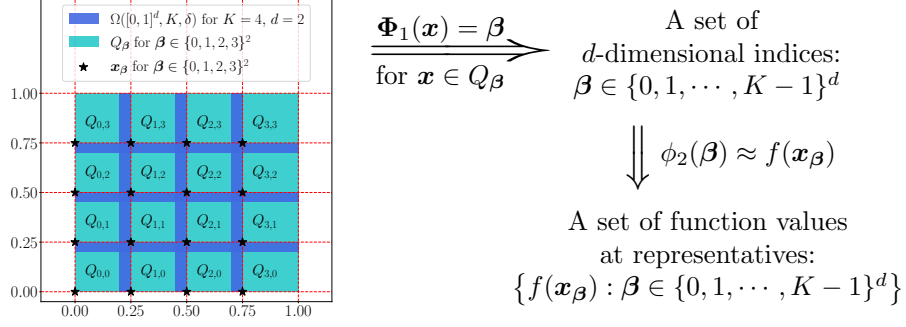
15

Figure 5: An illustration of the desired function $\phi = \phi_2 \circ \boldsymbol{\Phi}_1$. Note that $\phi \approx f$ on $[0,1]^d \backslash \Omega([0,1]^d, K, \delta)$, since $\phi(\boldsymbol{x}) = \phi_2 \circ \boldsymbol{\Phi}_1(\boldsymbol{x}) = \phi_2(\boldsymbol{\beta}) \approx f(\boldsymbol{x}_{\boldsymbol{\beta}}) \approx f(\boldsymbol{x})$ for any $\boldsymbol{x} \in Q_{\boldsymbol{\beta}}$ and each $\boldsymbol{\beta} \in \{0, 1, \cdots, K-1\}^d$.

**Proposition 3.1.** *For any $N, L, d \in \mathbb{N}^+$ and $\delta \in (0, \frac{1}{3K}]$ with*

$$K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor \lfloor n^{1/d} \rfloor, \quad \text{where } n = \lfloor \log_3(N+2) \rfloor,$$

*there exists a one-dimensional function $\phi$ implemented by a ReLU network with width $8 \lfloor N^{1/d} \rfloor + 3$ and depth $2 \lfloor L^{1/d} \rfloor + 5$ such that*

$$\phi(x) = k, \quad \text{if } x \in \left[ \frac{k}{K}, \frac{k+1}{K} - \delta \cdot 1_{\{k \le K-2\}} \right] \text{ for } k = 0, 1, \cdots, K-1.$$

The setting $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor^2 \lfloor n^{1/d} \rfloor = \mathcal{O}(N^{2/d} L^{2/d} n^{1/d})$ is not neat here, but it is very convenient for later use. The construction of $\phi_2$ is a direct result of Proposition 3.2 below, the proof of which relies on the bit extraction technique in [3].

**Proposition 3.2.** *Given any $\varepsilon > 0$ and arbitrary $N, L, J \in \mathbb{N}^+$ with $J \le N^2 L^2 \lfloor \log_3(N+2) \rfloor$, assume $y_j \ge 0$ for $j = 0, 1, \cdots, J-1$ are samples with*

$$|y_j - y_{j-1}| \le \varepsilon, \quad \text{for } j = 1, 2, \cdots, J-1.$$

*Then there exists $\phi \in \mathcal{NN}(\#\text{input} = 1; \text{width} \le 16N + 30; \text{depth} \le 6L + 10; \#\text{output} = 1)$ such that*

*(i) $|\phi(j) - y_j| \le \varepsilon$ for $j = 0, 1, \cdots, J-1$.*

*(ii) $0 \le \phi(x) \le \max\{y_j : j = 0, 1, \cdots, J-1\}$ for any $x \in \mathbb{R}$.*

With the above propositions ready, let us prove Theorem 2.1 in Section 3.2.

## 3.2 Constructive proof

We essentially construct an almost piecewise constant function implemented by a ReLU network with width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ to approximate $f$. We may assume $f$ is not a constant function since it is a trivial case. Then $\omega_f(r) > 0$ for any $r > 0$. It is clear that $|f(\boldsymbol{x}) - f(\boldsymbol{0})| \le \omega_f(\sqrt{d})$ for any $\boldsymbol{x} \in [0,1]^d$. Define $\widetilde{f} = f - f(\boldsymbol{0}) + \omega_f(\sqrt{d})$, then $0 \le \widetilde{f}(\boldsymbol{x}) \le 2\omega_f(\sqrt{d})$ for any $\boldsymbol{x} \in [0,1]^d$.

Let $M = N^2 L$, $n = \lfloor \log_3(N+2) \rfloor$, $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor^2 \lfloor n^{1/d} \rfloor$, and $\delta$ be an arbitrary number in $(0, \frac{1}{3K}]$. The proof can be divided into four steps as follows:

16

1. Normalize $f$ as $\widetilde{f}$, divide $[0,1]^d$ into a union of sub-cubes $\{Q_{\boldsymbol{\beta}}\}_{\boldsymbol{\beta} \in \{0,1,\cdots,K-1\}^d}$ and the trifling region $\Omega([0,1]^d, K, \delta)$, and denote $\boldsymbol{x}_{\boldsymbol{\beta}}$ as the vertex of $Q_{\boldsymbol{\beta}}$ with minimum $\|\cdot\|_1$ norm;

2. Construct a sub-network to implement a vector function $\boldsymbol{\Phi}_1$ projecting the whole cube $Q_{\boldsymbol{\beta}}$ to the $d$-dimensional index $\boldsymbol{\beta}$ for each $\boldsymbol{\beta}$, i.e., $\boldsymbol{\Phi}_1(\boldsymbol{x}) = \boldsymbol{\beta}$ for all $\boldsymbol{x} \in Q_{\boldsymbol{\beta}}$;

3. Construct a sub-network to implement a function $\phi_2$ mapping the index $\boldsymbol{\beta}$ approximately to $\widetilde{f}(\boldsymbol{x}_{\boldsymbol{\beta}})$. This core step can be further divided into three sub-steps:

    3.1. Construct a sub-network to implement $\psi_1$ bijectively mapping the index set $\{0, 1, \cdots, K-1\}^d$ to an auxiliary set $\mathcal{A}_1 \subseteq \left\{ \frac{j}{2K^d} : j = 0, 1, \cdots, 2K^d \right\}$ defined later (see Figure 7 for an illustration);

    3.2. Determine a continuous piecewise linear function $g$ with a set of breakpoints $\mathcal{A}_1 \cup \mathcal{A}_2 \cup \{1\}$ satisfying: 1) assign the values of $g$ at breakpoints in $\mathcal{A}_1$ based on $\{\widetilde{f}(\boldsymbol{x}_{\boldsymbol{\beta}})\}_{\boldsymbol{\beta}}$, i.e., $g \circ \psi_1(\boldsymbol{\beta}) = \widetilde{f}(\boldsymbol{x}_{\boldsymbol{\beta}})$; 2) assign the values of $g$ at breakpoints in $\mathcal{A}_2 \cup \{1\}$ to reduce the variation of $g$ for applying Proposition 3.2;

    3.3. Apply Proposition 3.2 to construct a sub-network to implement a function $\psi_2$ approximating $g$ well on $\mathcal{A}_1 \cup \mathcal{A}_2 \cup \{1\}$. Then the desired function $\phi_2$ is given by $\phi_2 = \psi_2 \circ \psi_1$ satisfying $\phi_2(\boldsymbol{\beta}) = \psi_2 \circ \psi_1(\boldsymbol{\beta}) \approx g \circ \psi_1(\boldsymbol{\beta}) = \widetilde{f}(\boldsymbol{x}_{\boldsymbol{\beta}})$;

4. Construct the final network to implement the desired function $\phi$ such that $\phi(\boldsymbol{x}) = \phi_2 \circ \boldsymbol{\Phi}_1(\boldsymbol{x}) + f(\boldsymbol{0}) - \omega_f(\sqrt{d}) \approx \widetilde{f}(\boldsymbol{x}_{\boldsymbol{\beta}}) + f(\boldsymbol{0}) - \omega_f(\sqrt{d}) = f(\boldsymbol{x}_{\boldsymbol{\beta}}) \approx f(\boldsymbol{x})$ for any $\boldsymbol{x} \in Q_{\boldsymbol{\beta}}$ and $\boldsymbol{\beta} \in \{0, 1, \cdots, K-1\}^d$.

The details of these steps can be found below.

**Step** 1: Divide $[0,1]^d$ into $\{Q_{\boldsymbol{\beta}}\}_{\boldsymbol{\beta} \in \{0,1,\cdots,K-1\}^d}$ and $\Omega([0,1]^d, K, \delta)$.

Define $\boldsymbol{x}_{\boldsymbol{\beta}} \coloneqq \boldsymbol{\beta}/K$ and

$$Q_{\boldsymbol{\beta}} \coloneqq \left\{ \boldsymbol{x} = [x_1, \cdots, x_d]^T \in [0,1]^d : x_i \in \left[ \tfrac{\beta_i}{K}, \tfrac{\beta_i+1}{K} - \delta \cdot 1_{\{\beta_i \le K-2\}} \right], \ i = 1, \cdots, d \right\}$$

for each $d$-dimensional index $\boldsymbol{\beta} = [\beta_1, \cdots, \beta_d]^T \in \{0, 1, \cdots, K-1\}^d$. Recall that $\Omega([0,1]^d, K, \delta)$ is the trifling region defined in Equation (2.1). Apparently, $\boldsymbol{x}_{\boldsymbol{\beta}}$ is the vertex of $Q_{\boldsymbol{\beta}}$ with minimum $\|\cdot\|_1$ norm and

$$[0,1]^d = \left( \cup_{\boldsymbol{\beta} \in \{0,1,\cdots,K-1\}^d} Q_{\boldsymbol{\beta}} \right) \bigcup \Omega([0,1]^d, K, \delta),$$

see Figure 6 for illustrations.

**Step** 2: Construct $\boldsymbol{\Phi}_1$ mapping $\boldsymbol{x} \in Q_{\boldsymbol{\beta}}$ to $\boldsymbol{\beta}$.

By Proposition 3.1, there exists $\phi_1 \in \mathcal{NN}(\text{width} \le 8\lfloor N^{1/d} \rfloor + 3; \text{depth} \le 2\lfloor L^{1/d} \rfloor + 5)$ such that

$$\phi_1(x) = k, \quad \text{if } x \in \left[ \tfrac{k}{K}, \tfrac{k+1}{K} - \delta \cdot 1_{\{k \le K-2\}} \right] \text{ for } k = 0, 1, \cdots, K-1.$$

It follows that $\phi_1(x_i) = \beta_i$ if $\boldsymbol{x} = [x_1, x_2, \cdots, x_d]^T \in Q_{\boldsymbol{\beta}}$ for each $\boldsymbol{\beta} = [\beta_1, \beta_2, \cdots, \beta_d]^T$.

17

Figure 6: Illustrations of $\Omega([0,1]^d, K, \delta)$, $Q_{\boldsymbol{\beta}}$, and $\boldsymbol{x}_{\boldsymbol{\beta}}$ for $\boldsymbol{\beta} \in \{0, 1, \cdots, K-1\}^d$. (a) $K = 4$ and $d = 1$. (b) $K = 4$ and $d = 2$.

By defining

$$\boldsymbol{\Phi}_1(\boldsymbol{x}) \coloneqq \big[\phi_1(x_1), \phi_1(x_2), \cdots, \phi_1(x_d)\big]^T, \quad \text{for any } \boldsymbol{x} = [x_1, x_2, \cdots, x_d]^T \in \mathbb{R}^d,$$

we have $\boldsymbol{\Phi}_1(\boldsymbol{x}) = \boldsymbol{\beta}$ if $\boldsymbol{x} \in Q_{\boldsymbol{\beta}}$ for $\boldsymbol{\beta} \in \{0, 1, \cdots, K-1\}^d$.

**Step** 3: Construct $\phi_2$ mapping $\boldsymbol{\beta}$ approximately to $\widetilde{f}(\boldsymbol{x}_{\boldsymbol{\beta}})$.

The construction of the sub-network implementing $\phi_2$ is essentially based on Proposition 3.2. To meet the requirements of applying Proposition 3.2, we first define two auxiliary set $\mathcal{A}_1$ and $\mathcal{A}_2$ as

$$\mathcal{A}_1 \coloneqq \Big\{ \tfrac{i}{K^{d-1}} + \tfrac{k}{2K^d} : i = 0, 1, \cdots, K^{d-1}-1 \quad \text{and} \quad k = 0, 1, \cdots, K-1 \Big\}$$

and

$$\mathcal{A}_2 \coloneqq \Big\{ \tfrac{i}{K^{d-1}} + \tfrac{K+k}{2K^d} : i = 0, 1, \cdots, K^{d-1}-1 \quad \text{and} \quad k = 0, 1, \cdots, K-1 \Big\}.$$

Clearly, $\mathcal{A}_1 \cup \mathcal{A}_2 \cup \{1\} = \big\{ \tfrac{j}{2K^d} : j = 0, 1, \cdots, 2K^d \big\}$ and $\mathcal{A}_1 \cap \mathcal{A}_2 = \varnothing$. See Figure 6 for an illustration of $\mathcal{A}_1$ and $\mathcal{A}_2$. Next, we further divide this step into three sub-steps.

**Step** 3.1: Construct $\psi_1$ bijectively mapping $\{0, 1, \cdots, K-1\}^d$ to $\mathcal{A}_1$.

Inspired by the binary representation, we define

$$\psi_1(\boldsymbol{x}) \coloneqq \frac{x_d}{2K^d} + \sum_{i=1}^{d-1} \frac{x_i}{K^i}, \quad \text{for any } \boldsymbol{x} = [x_1, x_2, \cdots, x_d]^T \in \mathbb{R}^d. \tag{3.1}$$

Then $\psi_1$ is a linear function bijectively mapping the index set $\{0, 1, \cdots, K-1\}^d$ to

$$\Big\{ \tfrac{\beta_d}{2K^d} + \sum_{i=1}^{d-1} \tfrac{\beta_i}{K^i} : \boldsymbol{\beta} \in \{0, 1, \cdots, K-1\}^d \Big\}$$

$$= \Big\{ \tfrac{i}{K^{d-1}} + \tfrac{k}{2K^d} : i = 0, 1, \cdots, K^{d-1}-1 \quad \text{and} \quad k = 0, 1, \cdots, K-1 \Big\} = \mathcal{A}_1.$$

**Step** 3.2: Construct $g$ to satisfy $g \circ \psi_1(\boldsymbol{\beta}) = \widetilde{f}(\boldsymbol{x}_{\boldsymbol{\beta}})$ and to meet the requirements of applying Proposition 3.2.
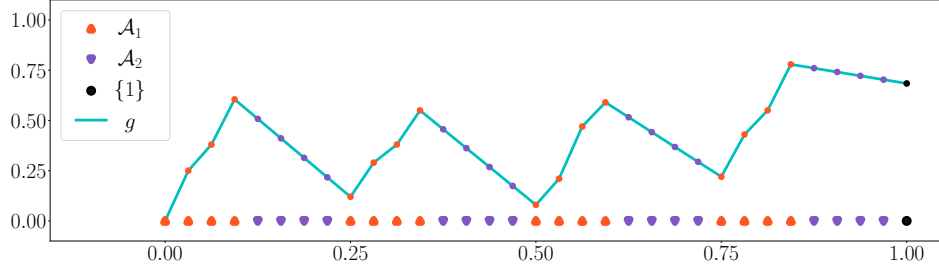
18

Figure 7: An illustration of $\mathcal{A}_1$, $\mathcal{A}_2$, $\{1\}$, and $g$ for $d = 2$ and $K = 4$.

Let $g : [0, 1] \to \mathbb{R}$ be a continuous piecewise linear function with a set of breakpoints $\left\{ \frac{j}{2K^d} : j = 0, 1, \cdots, 2K^d \right\} = \mathcal{A}_1 \cup \mathcal{A}_2 \cup \{1\}$ and the values of $g$ at these breakpoints satisfy the following properties:

- The values of $g$ at the breakpoints in $\mathcal{A}_1$ are set as

$$g(\psi_1(\boldsymbol{\beta})) = \widetilde{f}(\boldsymbol{x}_{\boldsymbol{\beta}}), \quad \text{for any } \boldsymbol{\beta} \in \{0, 1, \cdots, K-1\}^d; \tag{3.2}$$

- At the breakpoint 1, let $g(1) = \widetilde{f}(\mathbf{1})$, where $\mathbf{1} = [1, 1, \cdots, 1]^T \in \mathbb{R}^d$;

- The values of $g$ at the breakpoints in $\mathcal{A}_2$ are assigned to reduce the variation of $g$, which is a requirement of applying Proposition 3.2. Note that

$$\left\{ \frac{i}{K^{d-1}} - \frac{K+1}{2K^d}, \ \frac{i}{K^{d-1}} \right\} \subseteq \mathcal{A}_1 \cup \{1\}, \quad \text{for } i = 1, 2, \cdots, K^{d-1},$$

  implying the values of $g$ at $\frac{i}{K^{d-1}} - \frac{K+1}{2K^d}$ and $\frac{i}{K^{d-1}}$ have been assigned for $i = 1, 2, \cdots, K^{d-1}$. Thus, the values of $g$ at the breakpoints in $\mathcal{A}_2$ can be successfully assigned by letting $g$ linear on each interval $\left[ \frac{i}{K^{d-1}} - \frac{K+1}{2K^d}, \frac{i}{K^{d-1}} \right]$ for $i = 1, 2, \cdots, K^{d-1}$, since $\mathcal{A}_2 \subseteq \cup_{i=1}^{K^{d-1}} \left[ \frac{i}{K^{d-1}} - \frac{K+1}{2K^d}, \frac{i}{K^{d-1}} \right]$.

Apparently, such a function $g$ exists (see Figure 7 for an example) and satisfies

$$\left| g\left( \tfrac{j}{2K^d} \right) - g\left( \tfrac{j-1}{2K^d} \right) \right| \leq \max \left\{ \omega_f\left( \tfrac{1}{K} \right), \omega_f(\sqrt{d})/K \right\} \leq \omega_f\left( \tfrac{\sqrt{d}}{K} \right), \quad \text{for } j = 1, 2, \cdots, 2K^d,$$

and

$$0 \leq g\left( \tfrac{j}{2K^d} \right) \leq 2\omega_f(\sqrt{d}), \quad \text{for } j = 0, 1, \cdots, 2K^d.$$

**Step** 3.3: Construct $\psi_2$ approximating $g$ well on $\mathcal{A}_1 \cup \mathcal{A}_2 \cup \{1\}$.

Note that

$$2K^d = 2\left( \lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor^2 \lfloor n^{1/d} \rfloor \right)^d \leq 2\left( N^2 L^2 n \right) \leq N^2 \lceil \sqrt{2}L \rceil^2 \lfloor \log_3(N+2) \rfloor.$$

By Proposition 3.2 (set $y_j = g\left( \tfrac{j}{2K^2} \right)$ and $\varepsilon = \omega_f\left( \tfrac{\sqrt{d}}{K} \right) > 0$ therein), there exists

$$\widetilde{\psi}_2 \in \mathcal{NN}(\#\text{input} = 1; \ \text{width} \leq 16N + 30; \ \text{depth} \leq 6\lceil \sqrt{2}L \rceil + 10; \ \#\text{output} = 1)$$

19

such that

$$|\widetilde{\psi}_2(j) - g(\tfrac{j}{2K^d})| \le \omega_f(\tfrac{\sqrt{d}}{K}), \quad \text{for } j = 0, 1, \cdots, 2K^d - 1,$$

and

$$0 \le \widetilde{\psi}_2(x) \le \max\{g(\tfrac{j}{2K^d}) : j = 0, 1, \cdots, 2K^d - 1\} \le 2\omega_f(\sqrt{d}), \quad \text{for any } x \in \mathbb{R}.$$

By defining $\psi_2(x) \coloneqq \widetilde{\psi}_2(2K^d x)$ for any $x \in \mathbb{R}$, we have $\psi_2 \in \mathcal{NN}(\#\text{input} = 1; \text{ width} \le 16N + 30; \text{ depth} \le 6\lceil \sqrt{2}L \rceil + 10; \#\text{output} = 1)$,

$$0 \le \psi_2(x) = \widetilde{\psi}_2(2K^d x) \le 2\omega_f(\sqrt{d}), \quad \text{for any } x \in \mathbb{R}, \tag{3.3}$$

and

$$|\psi_2(\tfrac{j}{2K^d}) - g(\tfrac{j}{2K^d})| = |\widetilde{\psi}_2(j) - g(\tfrac{j}{2K^d})| \le \omega_f(\tfrac{\sqrt{d}}{K}), \quad \text{for } j = 0, 1, \cdots, 2K^d - 1. \tag{3.4}$$

Let us end Step 3 by defining the desired function $\phi_2$ as $\phi_2 \coloneqq \psi_2 \circ \psi_1$. Note that $\psi_1 : \mathbb{R}^d \to \mathbb{R}$ is a linear function and $\psi_2 \in \mathcal{NN}(\#\text{input} = 1; \text{ width} \le 16N + 30; \text{ depth} \le 6\lceil \sqrt{2}L \rceil + 10; \#\text{output} = 1)$. Thus, $\phi_2 \in \mathcal{NN}(\#\text{input} = 1; \text{ width} \le 16N + 30; \text{ depth} \le 6\lceil \sqrt{2}L \rceil + 10; \#\text{output} = 1)$. By Equation (3.2) and (3.4), we have

$$|\phi_2(\boldsymbol{\beta}) - \widetilde{f}(\boldsymbol{x_\beta})| = |\psi_2(\psi_1(\boldsymbol{\beta})) - g(\psi_1(\boldsymbol{\beta}))| \le \omega_f(\tfrac{\sqrt{d}}{K}), \tag{3.5}$$

for any $\boldsymbol{\beta} \in \{0, 1, \cdots, K - 1\}^d$. Equation (3.3) and $\phi_2 = \psi_2 \circ \psi_1$ implies

$$0 \le \phi_2(\boldsymbol{x}) \le 2\omega_f(\sqrt{d}), \quad \text{for any } \boldsymbol{x} \in \mathbb{R}^d. \tag{3.6}$$

**Step** 4: Construct the final network to implement the desired function $\phi$.

Define $\phi \coloneqq \phi_2 \circ \boldsymbol{\Phi}_1 + f(\boldsymbol{0}) - \omega_f(\sqrt{d})$. Since $\phi_1 \in \mathcal{NN}(\text{width} \le 8\lfloor N^{1/d} \rfloor + 3; \text{ depth} \le 2\lfloor L^{1/d} \rfloor + 5)$, we have $\boldsymbol{\Phi}_1 \in \mathcal{NN}(\#\text{input} = d; \text{ width} \le 8d\lfloor N^{1/d} \rfloor + 3d; \text{ depth} \le 2L + 5; \#\text{output} = d)$. If follows from the fact $\lceil \sqrt{2}L \rceil \le \lceil \tfrac{3}{2}L \rceil \le \tfrac{3}{2}L + \tfrac{1}{2}$ that $6\lceil \sqrt{2}L \rceil + 10 \le 9L + 13$, implying

$$\phi_2 \in \mathcal{NN}(\#\text{input} = 1; \text{ width} \le 16N + 30; \text{ depth} \le 6\lceil \sqrt{2}L \rceil + 10; \#\text{output} = 1)$$
$$\subseteq \mathcal{NN}(\#\text{input} = 1; \text{ width} \le 16N + 30; \text{ depth} \le 9L + 13; \#\text{output} = 1).$$

Thus, $\phi = \phi_2 \circ \boldsymbol{\Phi}_1 + f(\boldsymbol{0}) - \omega_f(\sqrt{d})$ is in

$$\mathcal{NN}\big(\text{width} \le \max\{8d\lfloor N^{1/d} \rfloor + 3d, 16N + 30\}; \text{ depth} \le (2L + 5) + (9L + 13) = 11L + 18\big).$$

Now let us estimate the approximation error. Note that $f = \widetilde{f} + f(\boldsymbol{0}) - \omega_f(\sqrt{d})$. By Equation (3.5), for any $\boldsymbol{x} \in Q_{\boldsymbol{\beta}}$ and $\boldsymbol{\beta} \in \{0, 1, \cdots, K - 1\}^d$, we have

$$|f(\boldsymbol{x}) - \phi(\boldsymbol{x})| = |\widetilde{f}(\boldsymbol{x}) - \phi_2(\boldsymbol{\Phi}_1(\boldsymbol{x}))| = |\widetilde{f}(\boldsymbol{x}) - \phi_2(\boldsymbol{\beta})|$$
$$\le |\widetilde{f}(\boldsymbol{x}) - \widetilde{f}(\boldsymbol{x_\beta})| + |\widetilde{f}(\boldsymbol{x_\beta}) - \phi_2(\boldsymbol{\beta})|$$
$$\le \omega_f(\tfrac{\sqrt{d}}{K}) + \omega_f(\tfrac{\sqrt{d}}{K}) \le 2\omega_f\big(64\sqrt{d}\big(N^2 L^2 \log_3(N + 2)\big)^{-1/d}\big),$$

20

where the last inequality comes from the fact

$$K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor^2 \lfloor n^{1/d} \rfloor \geq \frac{N^{2/d} L^{2/d} n^{1/d}}{32} = \frac{N^{2/d} L^{2/d} \lfloor \log_3(N+2) \rfloor^{1/d}}{32} \geq \frac{(N^2 L^2 \log_3(N+2))^{1/d}}{64},$$

for any $N, L \in \mathbb{N}^+$. Recall the fact $\omega_f(j \cdot r) \leq j \cdot \omega_f(r)$ for any $j \in \mathbb{N}^+$ and $r \in [0, \infty)$. Therefore, for any $\boldsymbol{x} \in \bigcup_{\boldsymbol{\beta} \in \{0,1,\cdots,K-1\}^d} Q_{\boldsymbol{\beta}} = [0,1]^d \backslash \Omega([0,1]^d, K, \delta)$, we have

$$
\begin{aligned}
|f(\boldsymbol{x}) - \phi(\boldsymbol{x})| &\leq 2\omega_f\left( 64\sqrt{d} \left( N^2 L^2 \log_3(N+2) \right)^{-1/d} \right) \\
&\leq 2\lceil 64\sqrt{d} \rceil \omega_f\left( \left( N^2 L^2 \log_3(N+2) \right)^{-1/d} \right) \\
&\leq 130\sqrt{d}\, \omega_f\left( \left( N^2 L^2 \log_3(N+2) \right)^{-1/d} \right).
\end{aligned}
$$

It remains to show the upper bound of $\phi$. By Equation (3.6) and $\phi = \phi_2 \circ \boldsymbol{\Phi}_1 + f(\boldsymbol{0}) - \omega_f(\sqrt{d})$, it holds that $\|\phi\|_{L^\infty(\mathbb{R}^d)} \leq |f(\boldsymbol{0})| + \omega_f(\sqrt{d})$. Thus, we finish the proof.

# 4 Proofs of propositions in Section 3.1

In this section, we will prove Proposition 3.1 and 3.2. We first introduce several basic results of ReLU networks. Next, we prove these two propositions based on these basic results.

## 4.1 Basic results of ReLU networks

To simplify the proofs of two propositions in Section 3.1, we introduce three lemmas below, which are basic results of ReLU networks

**Lemma 4.1.** *For any* $N_1, N_2 \in \mathbb{N}^+$, *given* $N_1(N_2 + 1) + 1$ *samples* $(x_i, y_i) \in \mathbb{R}^2$ *with* $x_0 < x_1 < \cdots < x_{N_1(N_2+1)}$ *and* $y_i \geq 0$ *for* $i = 0, 1, \cdots, N_1(N_2+1)$, *there exists* $\phi \in \mathcal{NN}(\#\text{input} = 1; \text{widthvec} = [2N_1, 2N_2 + 1]; \#\text{output} = 1)$ *satisfying the following conditions.*

*(i)* $\phi(x_i) = y_i$ *for* $i = 0, 1, \cdots, N_1(N_2 + 1)$.

*(ii)* $\phi$ *is linear on each interval* $[x_{i-1}, x_i]$ *for* $i \notin \{(N_2 + 1)j : j = 1, 2, \cdots, N_1\}$.

**Lemma 4.2.** *Given any* $N, L, d \in \mathbb{N}^+$, *it holds that*

$$
\begin{aligned}
&\mathcal{NN}(\#\text{input} = d; \text{widthvec} = [N, NL]; \#\text{output} = 1) \\
&\subseteq \mathcal{NN}(\#\text{input} = d; \text{width} \leq 2N + 2; \text{depth} \leq L + 1; \#\text{output} = 1).
\end{aligned}
$$

**Lemma 4.3.** *For any* $n \in \mathbb{N}^+$, *it holds that*

$$\text{CPwL}(\mathbb{R}, n) \subseteq \mathcal{NN}(\#\text{input} = 1; \text{widthvec} = [n+1]; \#\text{output} = 1). \tag{4.1}$$

Lemma 4.1 is a part of Theorem 3.2 in [29] or Lemma 2.2 in [22]. Lemma 4.1 is Theorem 3.1 in [29] or Lemma 3.4 in [22]. It remains to prove Lemma 4.3.

*Proof of Lemma 4.3.* We use the mathematics induction to prove Equation (4.1). First, consider the case $n = 1$. Given any $f \in \mathrm{CPwL}(\mathbb{R}, n)$, there exist $a_1, a_2, x_0 \in \mathbb{R}$ such that

$$f(x) = \begin{cases} a_1(x - x_0) + f(x_0), & \text{if } x \geq x_0, \\ a_2(x_0 - x) + f(x_0), & \text{if } x < x_0. \end{cases}$$

Thus, $f(x) = a_1\sigma(x - x_0) + a_2\sigma(x_0 - x) + f(x_0)$ for any $x \in \mathbb{R}$, implying

$$f \in \mathcal{NN}(\#\text{input} = 1; \text{ widthvec} = [2]; \#\text{output} = 1).$$

Thus, Equation (4.1) holds for $n = 1$.

Now assume Equation (4.1) holds for $n = k \in \mathbb{N}^+$, we would like to show it is also true for $n = k + 1$. Given any $f \in \mathrm{CPwL}(\mathbb{R}, k + 1)$, we may assume the biggest breakpoint of $f$ is $x_0$ since it is trivial for the case that $f$ has no breakpoint. Denote the slopes of the linear pieces left and right next to $x_0$ by $a_1$ and $a_2$, respectively. Define

$$\widetilde{f}(x) := f(x) - (a_2 - a_1)\sigma(x - x_0), \quad \text{for any } x \in \mathbb{R}.$$

Then $\widetilde{f}$ has at most $k$ breakpoints. By the induction hypothesis, we have

$$\widetilde{f} \in \mathrm{CPwL}(\mathbb{R}, k) \subseteq \mathcal{NN}(\#\text{input} = 1; \text{ widthvec} = [k + 1]; \#\text{output} = 1).$$

Thus, there exist $w_{0,j}, b_{0,j}, w_{1,j}, b_1$ for $j = 1, 2, \cdots, k + 1$ such that

$$\widetilde{f}(x) = \sum_{j=1}^{k+1} w_{1,j}\sigma(w_{0,j}x + b_{0,j}) + b_1, \quad \text{for any } x \in \mathbb{R}.$$

Therefore, for any $x \in \mathbb{R}$, we have

$$f(x) = (a_2 - a_1)\sigma(x - x_0) + \widetilde{f}(x) = (a_2 - a_1)\sigma(x - x_0) + \sum_{j=1}^{k} w_{1,j}\sigma(w_{0,j}x + b_{0,j}) + b_1,$$

implying $f \in \mathcal{NN}(\#\text{input} = 1; \text{ widthvec} = [k + 1]; \#\text{output} = 1)$. Thus, Equation (4.1) holds for $k+1$, which means we finish the induction process. So we complete the proof. $\square$

## 4.2 Proof of Proposition 3.1

Now, let us present the detailed proof of Proposition 3.1. Denote $K = \widetilde{M} \cdot \widetilde{L}$, where $\widetilde{M} = \lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor$, $n = \lfloor \log_3(N + 2) \rfloor$, and $\widetilde{L} = \lfloor L^{1/d} \rfloor \lfloor n^{1/d} \rfloor$. Consider the sample set

$$\left\{ (1, \widetilde{M} - 1), (2, 0) \right\} \bigcup \left\{ (\tfrac{m}{\widetilde{M}}, m) : m = 0, 1, \cdots, \widetilde{M} - 1 \right\}$$
$$\bigcup \left\{ (\tfrac{m+1}{\widetilde{M}} - \delta, m) : m = 0, 1, \cdots, \widetilde{M} - 2 \right\}.$$

Its size is

$$2\widetilde{M} + 1 = 2\lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor + 1 = \lfloor N^{1/d} \rfloor \cdot \left( \left( 2\lfloor N^{1/d} \rfloor \lfloor L^{1/d} \rfloor - 1 \right) + 1 \right) + 1.$$

By Lemma 4.1 (set $N_1 = \lfloor N^{1/d} \rfloor$ and $N_2 = 2\lfloor N^{1/d} \rfloor \lfloor L^{1/d} \rfloor - 1$ therein), there exists

$$\phi_1 \in \mathcal{NN}\left( \text{widthvec} = \left[ 2\lfloor N^{1/d} \rfloor, 2(2\lfloor N^{1/d} \rfloor \lfloor L^{1/d} \rfloor - 1) + 1 \right] \right)$$
$$= \mathcal{NN}\left( \text{widthvec} = \left[ 2\lfloor N^{1/d} \rfloor, 4\lfloor N^{1/d} \rfloor \lfloor L^{1/d} \rfloor - 1 \right] \right)$$

such that

- $\phi_1\left(\frac{\widetilde{M}-1}{\widetilde{M}}\right) = \phi_1(1) = \widetilde{M} - 1$ and $\phi_1\left(\frac{m}{\widetilde{M}}\right) = \phi_1\left(\frac{m+1}{\widetilde{M}} - \delta\right) = m$ for $m = 0, 1, \cdots, \widetilde{M} - 2$.

- $\phi_1$ is linear on $\left[\frac{\widetilde{M}-1}{\widetilde{M}}, 1\right]$ and each interval $\left[\frac{m}{\widetilde{M}}, \frac{m+1}{\widetilde{M}} - \delta\right]$ for $m = 0, 1, \cdots, \widetilde{M} - 2$.

Then, for $m = 0, 1, \cdots, \widetilde{M} - 1$, we have

$$\phi_1(x) = m, \quad \text{for any } x \in \left[\frac{m}{\widetilde{M}}, \frac{m+1}{\widetilde{M}} - \delta \cdot 1_{\{m \le \widetilde{M}-2\}}\right]. \tag{4.2}$$

Now consider the another sample set

$$\left\{\left(\frac{1}{\widetilde{M}}, \widetilde{L} - 1\right), (2, 0)\right\} \bigcup \left\{\left(\frac{\ell}{\widetilde{M}\widetilde{L}}, \ell\right) : \ell = 0, 1, \cdots, \widetilde{L} - 1\right\}$$
$$\bigcup \left\{\left(\frac{\ell+1}{\widetilde{M}\widetilde{L}} - \delta, \ell\right) : \ell = 0, 1, \cdots, \widetilde{L} - 2\right\}.$$

Its size is

$$2\widetilde{L} + 1 = 2\lfloor L^{1/d}\rfloor\lfloor n^{1/d}\rfloor + 1 = \lfloor n^{1/d}\rfloor \cdot \left(\left(2\lfloor L^{1/d}\rfloor - 1\right) + 1\right) + 1.$$

By Lemma 4.1 (set $N_1 = \lfloor n^{1/d}\rfloor$ and $N_2 = 2\lfloor L^{1/d}\rfloor - 1$ therein), there exists

$$\phi_2 \in \mathcal{NN}\left(\text{widthvec} = \left[2\lfloor n^{1/d}\rfloor, 2(2\lfloor L^{1/d}\rfloor - 1) + 1\right]\right)$$
$$= \mathcal{NN}\left(\text{widthvec} = \left[2\lfloor n^{1/d}\rfloor, 4\lfloor L^{1/d}\rfloor - 1\right]\right)$$

such that

- $\phi_2\left(\frac{\widetilde{L}-1}{\widetilde{M}\widetilde{L}}\right) = \phi_2\left(\frac{1}{\widetilde{M}}\right) = \widetilde{L} - 1$ and $\phi_2\left(\frac{\ell}{\widetilde{M}\widetilde{L}}\right) = \phi_2\left(\frac{\ell+1}{\widetilde{M}\widetilde{L}} - \delta\right) = \ell$ for $\ell = 0, 1, \cdots, \widetilde{L} - 2$.

- $\phi_2$ is linear on $\left[\frac{\widetilde{L}-1}{\widetilde{M}\widetilde{L}}, \frac{1}{\widetilde{M}}\right]$ and each interval $\left[\frac{\ell}{\widetilde{M}\widetilde{L}}, \frac{\ell+1}{\widetilde{M}\widetilde{L}} - \delta\right]$ for $\ell = 0, 1, \cdots, \widetilde{L} - 2$.

It follows that, for $m = 0, 1, \cdots, \widetilde{M} - 1$ and $\ell = 0, 1, \cdots, \widetilde{L} - 1$,

$$\phi_2\left(x - \frac{m}{\widetilde{M}}\right) = \ell, \quad \text{for any } x \in \left[\frac{m\widetilde{L}+\ell}{\widetilde{M}\widetilde{L}}, \frac{m\widetilde{L}+\ell+1}{\widetilde{M}\widetilde{L}} - \delta \cdot 1_{\{\ell \le \widetilde{L}-2\}}\right]. \tag{4.3}$$

$K = \widetilde{M} \cdot \widetilde{L}$ implies any $k \in \{0, 1, \cdots, K-1\}$ can be unique represented by $k = m\widetilde{L} + \ell$ for $m = 0, 1, \cdots, \widetilde{M} - 1$ and $\ell = 0, 1, \cdots, \widetilde{L} - 1$. Then the desired function $\phi$ can be implemented by a ReLU network shown in Figure 8.
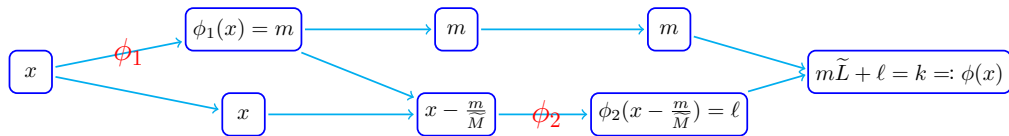


Figure 8: An illustration of the network architecture implementing $\phi$ based on Equation (4.2) and (4.3) for $x \in \left[\frac{k}{K}, \frac{k}{K} - \delta \cdot 1_{\{k \le K-2\}}\right] = \left[\frac{m\widetilde{L}+\ell}{\widetilde{M}\widetilde{L}}, \frac{m\widetilde{L}+\ell+1}{\widetilde{M}\widetilde{L}} - \delta \cdot 1_{\{m \le \widetilde{M}-2 \text{ or } \ell \le \widetilde{L}-2\}}\right]$, where $k = m\widetilde{L} + \ell$ for $m = 0, 1, \cdots, \widetilde{M} - 1$ and $\ell = 0, 1, \cdots, \widetilde{L} - 1$.

Clearly,

$$\phi(x) = k, \quad \text{if } x \in \left[\frac{k}{K}, \frac{k}{K} - \delta \cdot 1_{\{k \le K-2\}}\right], \quad \text{for any } k \in \{0, 1, \cdots, K - 1\}.$$

By Lemma 4.2, we have

$$
\begin{aligned}
\phi_1 \in \mathcal{NN}\big(&\#\text{input} = 1;\ \text{widthvec} = \big[2\lfloor N^{1/d}\rfloor, 4\lfloor N^{1/d}\rfloor\lfloor L^{1/d}\rfloor - 1\big];\ \#\text{output} = 1\big) \\
&\subseteq \mathcal{NN}\big(\#\text{input} = 1;\ \text{width} \le 8\lfloor N^{1/d}\rfloor + 2;\ \text{depth} \le \lfloor L^{1/d}\rfloor + 1;\ \#\text{output} = 1\big)
\end{aligned}
$$

and

$$
\begin{aligned}
\phi_2 \in \mathcal{NN}\big(&\#\text{input} = 1;\ \text{widthvec} = \big[2\lfloor n^{1/d}\rfloor, 4\lfloor L^{1/d}\rfloor - 1\big];\ \#\text{output} = 1\big) \\
&\subseteq \mathcal{NN}\big(\#\text{input} = 1;\ \text{width} \le 8\lfloor n^{1/d}\rfloor + 2;\ \text{depth} \le \lfloor L^{1/d}\rfloor + 1;\ \#\text{output} = 1\big).
\end{aligned}
$$

Recall that $n = \lfloor \log_3(N+2)\rfloor \le N$. It follows from Figure 8 that $\phi$ can be implemented by a ReLU network with width

$$
\max\big\{8\lfloor N^{1/d}\rfloor + 2 + 1, 8\lfloor n^{1/d}\rfloor + 2 + 1\big\} = 8\lfloor N^{1/d}\rfloor + 3
$$

and depth

$$
(\lfloor L^{1/d}\rfloor + 1) + 2 + (\lfloor L^{1/d}\rfloor + 1) + 1 = 2\lfloor L^{1/d}\rfloor + 5.
$$

So we finish the proof.

## 4.3 Proof of Proposition 3.2

The proof of Proposition 3.2 is based on the bit extraction technique in [3, 8]. In fact, we modify this technique to extract the sum of many bits rather than one bit and this modification can be summarized in Lemma 4.4 and 4.5 below.

**Lemma 4.4.** *For any $n \in \mathbb{N}^+$, there exists a function $\phi$ in*

$$
\mathcal{NN}\big(\#\text{input} = 2;\ \text{width} \le (n+1)2^{n+1};\ \text{depth} \le 3;\ \#\text{output} = 1\big)
$$

*such that: Given any $\theta_j \in \{0, 1\}$ for $j = 1, 2, \cdots, n$, we have*

$$
\phi(\text{bin}\,0.\theta_1\theta_2\cdots\theta_n, i) = \sum_{j=1}^{i} \theta_j, \quad \text{for any } i \in \{0, 1, 2, \cdots, n\}.^{②}
$$

*Proof.* Define $\theta = \text{bin}\,0.\theta_1\theta_2\cdots\theta_n$. Clearly,

$$
\theta_j = \lfloor 2^j\theta\rfloor/2 - \lfloor 2^{j-1}\theta\rfloor, \quad \text{for any } j \in \{1, 2, \cdots, n\}.
$$

We shall use a ReLU network to replace $\lfloor\cdot\rfloor$. Let $g \in \text{CPwL}(\mathbb{R}, 2^{n+1} - 2)$ be the function satisfying two conditions:

- $g$ matches set of samples

$$
\bigcup_{k=0}^{2^n - 1} \big\{(k, k), (k+1-\delta, k)\big\}, \quad \text{where } \delta = 2^{-(n+1)};
$$

---
② By convention, $\sum_{j=n}^{m} a_j = 0$ if $n > m$, no matter what $a_j$ is for each $j$.

- The breakpoint set of $g$ is

$$\Big( \bigcup_{k=0}^{2^n-1} \{k, k+1-\delta\} \Big) \Big\backslash \Big( \{0\} \cup \{2^n - \delta\} \Big).$$

Then $g(x) = \lfloor x \rfloor$ for any $x \in \bigcup_{k=0}^{2^n-1} [k, k+1-\delta]$. Clearly, $\theta = \mathrm{bin}\, 0.\theta_1\theta_2\cdots\theta_n$ implies

$$2^j \theta \in \bigcup_{k=0}^{2^n-1} [k, k+1-\delta], \quad \text{for any } j \in \{1, 2, \cdots, n\}.$$

Thus,

$$\theta_j = \lfloor 2^j \theta \rfloor / 2 - \lfloor 2^{j-1}\theta \rfloor = g(2^j\theta)/2 - g(2^{j-1}\theta), \quad \text{for any } j \in \{1, 2, \cdots, n\}. \tag{4.4}$$

It is easy to design a ReLU network to output $\theta_1, \theta_2, \cdots, \theta_n$ by Equation (4.4) when using $\theta = \mathrm{bin}\, 0.\theta_1\theta_2\cdots\theta_n$ as the input. However, it is highly non-trivial to construct a ReLU network to output $\sum_{j=1}^{i} \theta_j$ with another input $i$, since many operations like multiplication and comparison are not allowed in designing ReLU networks. Now let us establish a formula to represent $\sum_{j=1}^{i} \theta_j$ in a form of a ReLU network as follows.

Define $\mathcal{T}(n) := \sigma(n+1) - \sigma(n) = \begin{cases} 1, & n \geq 0, \\ 0, & n < 0, \end{cases}$ for any integer $n$. Then, by Equation (4.4) and the fact $x_1 x_2 = \sigma(x_1 + x_2 - 1)$ for any $x_1, x_2 \in \{0, 1\}$, we have, for $i = 0, 1, 2, \cdots, n$,

$$\begin{aligned}
\sum_{j=1}^{i} \theta_j &= \sum_{j=1}^{n} \theta_j \cdot \mathcal{T}(i-j) = \sum_{j=1}^{n} \theta_j \cdot \Big( \sigma(i-j+1) - \sigma(i-j) \Big) \\
&= \sum_{j=1}^{n} \sigma \Big( \theta_j + \sigma(i-j+1) - \sigma(i-j) - 1 \Big) \\
&= \sum_{j=1}^{n} \sigma \Big( g(2^j\theta)/2 - g(2^{j-1}\theta) + \sigma(i-j+1) - \sigma(i-j) - 1 \Big).
\end{aligned}$$

Define

$$z_{i,j} := \sigma \Big( g(2^j\theta)/2 - g(2^{j-1}\theta) + \sigma(i-j+1) - \sigma(i-j) - 1 \Big), \tag{4.5}$$

for any $i, j \in \{1, 2, \cdots, n\}$. Then the goal is to design $\phi$ satisfying

$$\phi(\theta, i) = \sum_{j=1}^{i} \theta_j = \sum_{j=1}^{n} z_{i,j}, \quad \text{for any } i \in \{0, 1, 2, \cdots, n\}. \tag{4.6}$$

See Figure 9 for the network architecture implementing the desired function $\phi$.

By Lemma 4.3, we have

$$g \in \mathrm{CPwL}\big(\mathbb{R}, 2^{n+1} - 2\big) \subseteq \mathcal{NN}\big(\#\mathrm{input} = 1;\ \mathrm{widthvec} = [2^{n+1} - 1];\ \#\mathrm{output} = 1\big),$$

implying

$$g(2^j \cdot) \in \mathrm{CPwL}\big(\mathbb{R}, 2^{n+1} - 2\big) \subseteq \mathcal{NN}\big(\#\mathrm{input} = 1;\ \mathrm{widthvec} = [2^{n+1} - 1];\ \#\mathrm{output} = 1\big),$$

for any $j = 0, 1, 2, \cdots, n$. Clearly, the network in Figure 9 has width $(n+1)(2^{n+1} - 1) + (n+1) = (n+1)2^{n+1}$ and depth 3. So we finish the proof. $\qquad\square$
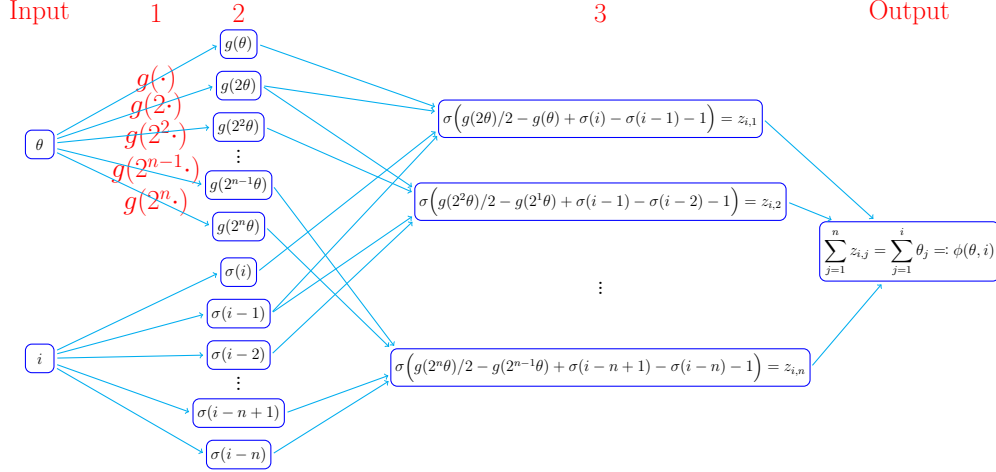
25

Figure 9: An illustration of the network implementing the desired function $\phi$ with the input $[\theta, i]^T = [\text{bin}\, 0.\theta_1\theta_2\cdots\theta_n, i]^T$ for any $i \in \{0, 1, 2, \cdots, n\}$ and $\theta_1, \theta_2, \cdots, \theta_n \in \{0, 1\}$. $g(2^j\cdot)$ can be implemented by a one-hidden-layer network with width $2^{n+1} - 1$ for each $j \in \{0, 1, \cdots, n\}$. The red numbers above the architecture indicate the order of hidden layers. The network architecture is essentially determined by Equation (4.5) and (4.6), which are valid no matter what $\theta_1, \theta_2, \cdots, \theta_n \in \{0, 1\}$ are. Thus, the desired function $\phi$ is independent of $\theta_1, \theta_2, \cdots, \theta_n \in \{0, 1\}$. We omit ReLU ($\sigma$) for a neuron if its output is non-negative without ReLU. Such a simplification are applied to similar figures in this paper.

**Lemma 4.5.** *For any $n, L \in \mathbb{N}^+$, there exists a function $\phi$ in*

$$\mathcal{NN}\big(\#\text{input} = 2; \ \text{width} \le (n + 3)2^{n+1} + 4; \ \text{depth} \le 4L + 2; \ \#\text{output} = 1\big)$$

*such that: Given any $\theta_j \in \{0, 1\}$ for $j = 1, 2, \cdots, Ln$, we have*

$$\phi(\text{bin}\, 0.\theta_1\theta_2\cdots\theta_{Ln}, k) = \sum_{j=1}^{k} \theta_j, \quad \text{for any } k \in \{1, 2, \cdots, Ln\}.$$

*Proof.* Let $g_1 \in \text{CPwL}(\mathbb{R}, 2^{n+1} - 2)$ be the function satisfying:

- $g_1$ matches the set of samples

$$\bigcup_{i=0}^{2^n - 1} \big\{(i, i), (i + 1 - \delta, i)\big\}, \quad \text{where } \delta = 2^{-(Ln+1)}.$$

- The breakpoint set of $g_1$ is

$$\Big( \bigcup_{i=0}^{2^n - 1} \big\{(i, i), (i + 1 - \delta, i)\big\}\Big) \Big\backslash \Big(\{0\} \cup \{2^n - \delta\}\Big).$$

Then $g_1(x) = \lfloor x \rfloor$ for any $x \in \cup_{i=0}^{2^n - 1}[i, i + 1 - \delta]$. Note that

$$2^n \cdot \text{bin}\, 0.\theta_{\ell n+1}\cdots\theta_{Ln} \in \bigcup_{i=0}^{2^n - 1}[i, i + 1 - \delta], \quad \text{for any } \ell \in \{0, 1, \cdots, L - 1\}.$$

Thus, for any $\ell \in \{0, 1, \cdots, L-1\}$, we have

$$\text{bin}\, 0.\theta_{\ell n+1}\cdots\theta_{\ell n+n} = \frac{\lfloor 2^n \cdot \text{bin}\, 0.\theta_{\ell n+1}\cdots\theta_{Ln}\rfloor}{2^n} = \frac{g_1(2^n \cdot \text{bin}\, 0.\theta_{\ell n+1}\cdots\theta_{Ln})}{2^n}. \tag{4.7}$$

Define $g_2(x) := 2^n x - g_1(2^n x)$ for any $x \in \mathbb{R}$. Then $g_2 \in \text{CPwL}(\mathbb{R}, 2^{n+1} - 2)$ and

$$\begin{aligned}
\text{bin}\, 0.\theta_{(\ell+1)n+1}\cdots\theta_{Ln} &= 2^n\Big(\text{bin}\, 0.\theta_{\ell n+1}\cdots\theta_{Ln} - \text{bin}\, 0.\theta_{\ell n+1}\cdots\theta_{\ell n+n}\Big) \\
&= 2^n\Big(\text{bin}\, 0.\theta_{\ell n+1}\cdots\theta_{Ln} - \frac{g_1(2^n \cdot \text{bin}\, 0.\theta_{\ell n+1}\cdots\theta_{Ln})}{2^n}\Big) = g_2(\text{bin}\, 0.\theta_{\ell n+1}\cdots\theta_{Ln}).
\end{aligned} \tag{4.8}$$

By Lemma 4.4, there exists

$$\phi_1 \in \mathcal{NN}\big(\#\text{input} = 2;\ \text{width} \leq (n+1)2^{n+1};\ \text{depth} \leq 3;\ \#\text{output} = 1\big)$$

such that: For any $\xi_1, \xi_2, \cdots, \xi_n \in \{0, 1\}$, we have

$$\phi_1\big(\text{bin}\, 0.\xi_1\xi_2\cdots\xi_n,\ i\big) = \sum_{j=1}^{i} \xi_j, \quad \text{for } i = 0, 1, 2, \cdots, n.$$

It follows that

$$\phi_1\big(\text{bin}\, 0.\theta_{\ell n+1}\theta_{\ell n+2}\cdots\theta_{\ell n+n},\ i\big) = \sum_{j=1}^{i} \theta_{\ell n+j}, \quad \text{for } \ell = 0, 1, \cdots, L-1 \text{ and } i = 0, 1, \cdots, n. \tag{4.9}$$

Define $\phi_{2,\ell}(x) := \min\{\sigma(x - \ell n), n\}$ for any $x \in \mathbb{R}$ and $\ell \in \{0, 1, \cdots, L-1\}$. For any $k \in \{1, 2, \cdots, Ln\}$, there exists $k_1 \in \{0, 1, \cdots, L-1\}$ and $k_2 \in \{1, 2, \cdots, n\}$ such that $k = k_1 n + k_2$, implying

$$\begin{aligned}
\sum_{i=1}^{k} \theta_i = \sum_{i=1}^{k_1 n + k_2} \theta_i &= \sum_{\ell=0}^{k_1-1}\Big(\sum_{j=1}^{n} \theta_{\ell n+j}\Big) + \sum_{\ell=k_1}^{k_1}\Big(\sum_{j=1}^{k_2} \theta_{\ell n+j}\Big) + \sum_{\ell=k_1+1}^{L-1}\Big(\sum_{j=1}^{0} \theta_{\ell n+j}\Big) \\
&= \sum_{\ell=0}^{L-1}\Big(\sum_{j=1}^{\min\{\sigma(k-\ell n),\, n\}} \theta_{\ell n+j}\Big) = \sum_{\ell=0}^{L-1}\Big(\sum_{j=1}^{\phi_{2,\ell}(k)} \theta_{\ell n+j}\Big).
\end{aligned} \tag{4.10}$$

Then, the desired function $\phi$ can be implemented by the network architecture in Figure 10.

By Lemma 4.3, we have

$$g_1, g_2 \in \text{CPwL}(\mathbb{R}, 2^{n+1} - 2) \subseteq \mathcal{NN}\big(\#\text{input} = 1;\ \text{widthvec} = [2^{n+1} - 1];\ \#\text{output} = 1\big).$$

Recall that $\phi_1 \in \mathcal{NN}\big(\text{width} \leq (n+1)2^{n+1};\ \text{depth} \leq 3\big)$. As shown in Figure 11, $\phi_{2,\ell}(x) \in \mathcal{NN}(\text{width} \leq 4;\ \text{depth} \leq 2)$ for $\ell = 0, 1, \cdots, L-1$. Therefore, the network in Figure 10 has width

$$(2^{n+1} - 1)\ +\ (2^{n+1} - 1)\ +\ (n+1)2^{n+1}\ +\ 1\ +\ 4\ +\ 1 = (n+3)2^{n+1} + 4$$

and depth

$$2 + L(1 + 3) = 4L + 2.$$
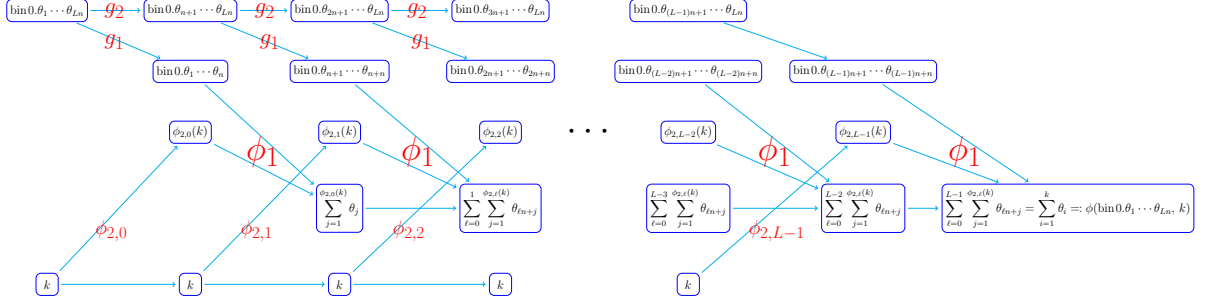
So we finish the proof. $\qquad\square$

Figure 10: An illustration of the network implementing the desired function $\phi$ with the input $[\text{bin}\,0.\theta_1\theta_2\cdots\theta_{Ln}, k]^T$ for any $k \in \{1, 2, \cdots, Ln\}$ and $\theta_1, \theta_2, \cdots, \theta_{Ln} \in \{0, 1\}$. The network architecture is essentially determined by Equation (4.7), (4.8), (4.9), and (4.10), which are valid no matter what $\theta_1, \theta_2, \cdots, \theta_{Ln} \in \{0, 1\}$ are. Thus, the desired function $\phi$ is independent of $\theta_1, \theta_2, \cdots, \theta_{Ln} \in \{0, 1\}$. We omit ReLU ($\sigma$) for a neuron if its output is non-negative without ReLU.



Figure 11: An illustration of the network implementing the desired function $\phi_{2,\ell}$ for each $\ell \in \{0, 1, \cdots, L-1\}$, based on $\min\{x, n\} = \frac{1}{2}\big(\sigma(x-n) - \sigma(-x-n) - \sigma(x-n) - \sigma(-x+n)\big)$.

Next, we introduce Lemma 4.6 to map indices to the partial sum of given bits.

**Lemma 4.6.** *Given any $N, L \in \mathbb{N}^+$ and arbitrary $\theta_{m,k} \in \{0, 1\}$ for $m = 0, 1, \cdots, M-1$ and $k = 0, 1, \cdots, Ln-1$, where $M = N^2L$ and $n = \lfloor \log_3(N+2) \rfloor$, there exists*

$$\phi \in \mathcal{NN}\big(\#\text{input} = 2; \ \text{width} \leq 6N + 14; \ \text{depth} \leq 5L + 4; \ \#\text{output} = 1\big)$$

*such that*

$$\phi(m, k) = \sum_{j=0}^{k} \theta_{m,j}, \quad \text{for } m = 0, 1, \cdots, M-1 \text{ and } k = 0, 1, \cdots, Ln-1.$$

*Proof.* Define

$$y_m := \text{bin}\,0.\theta_{m,0}\theta_{m,1}\cdots\theta_{m,Ln-1}, \quad \text{for } m = 0, 1, \cdots, M-1.$$

Consider the sample set $\{(m, y_m) : m = 0, 1, \cdots, M\}$, whose cardinality is

$$M + 1 = N\big((NL-1) + 1\big) + 1.$$

By Lemma 4.1 (set $N_1 = N$ and $N_2 = NL - 1$ therein), there exists

$$\phi_1 \in \mathcal{NN}\big(\#\text{input} = 1; \ \text{widthvec} = [2N, 2(NL-1) + 1]; \ \#\text{output} = 1\big)$$
$$= \mathcal{NN}\big(\#\text{input} = 1; \ \text{widthvec} = [2N, 2NL - 1]; \ \#\text{output} = 1\big)$$
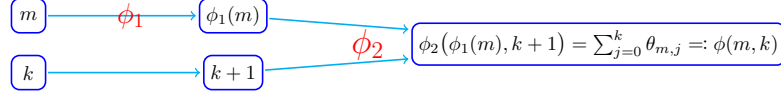
28

Figure 12: An illustration of the network implementing the desired function $\phi$ for $m = 0, 1, \cdots, M - 1$ and $k = 0, 1, \cdots, Ln - 1$.

such that

$$\phi_1(m) = y_m, \quad \text{for } m = 0, 1, \cdots, M - 1.$$

By Lemma 4.4, there exists

$$\phi_2 \in \mathcal{NN}\big(\#\text{input} = 2; \ \text{width} \leq (n + 3)2^{n+1} + 4; \ \text{depth} \leq 4L + 2; \ \#\text{output} = 1\big)$$

such that, for any $\xi_1, \xi_2, \cdots, \xi_{Ln} \in \{0, 1\}$, we have

$$\phi_2\big(\text{bin}\,0.\xi_1\xi_2\cdots\xi_{Ln}, \ k\big) = \sum_{j=1}^{k} \xi_j, \quad \text{for } k = 1, 2, \cdots, Ln.$$

It follows that, for any $\xi_0, \xi_1, \cdots, \xi_{Ln-1} \in \{0, 1\}$, we have

$$\phi_2\big(\text{bin}\,0.\xi_0\xi_1\cdots\xi_{Ln-1}, \ k + 1\big) = \sum_{j=0}^{k} \xi_j, \quad \text{for } k = 0, 1, \cdots, Ln - 1.$$

Thus, for $m = 0, 1, \cdots, M - 1$ and $k = 0, 1, \cdots, Ln - 1$, we have

$$\phi_2\big(\phi_1(m), k + 1\big) = \phi_2\big(y_m, k + 1\big) = \phi_2\big(0.\theta_{m,0}\theta_{m,1}\cdots\theta_{m,L-1}, k + 1\big) = \sum_{j=0}^{k} \theta_{m,j}.$$

Hence, the desired function function $\phi$ can be implemented by the network shown in Figure 12. By Lemma 4.2, $\phi_1 \in \mathcal{NN}(\text{widthvec} = [2N, 2NL - 1]) \subseteq \mathcal{NN}(\text{width} \leq 4N + 2; \ \text{depth} \leq L + 1)$. It holds that

$$(n + 3)2^{n+1} + 4 \leq 6 \cdot (3^n) + 2 = 6 \cdot \big(3^{\lfloor \log_3(N+2) \rfloor}\big) + 2 \leq 6(N + 2) + 2 = 6N + 14,$$

implying

$$\phi_2 \in \mathcal{NN}\big(\#\text{input} = 2; \ \text{width} \leq (n + 3)2^{n+1} + 4; \ \text{depth} \leq 4L + 2; \ \#\text{output} = 1\big)$$
$$\subseteq \mathcal{NN}\big(\#\text{input} = 2; \ \text{width} \leq 6N + 14; \ \text{depth} \leq 4L + 2; \ \#\text{output} = 1\big).$$

Therefore, the network in Figure 12 is with width $\max\{(4N + 2) + 1, 6N + 14\} = 6N + 14$ and depth $(4L + 2) + 1 + (L + 1) = 5L + 4$. So we finish the proof. $\quad\square$

Next, we apply Lemma 4.6 to prove Lemma 4.7 below, which is a key intermediate conclusion to prove Proposition 3.2.

**Lemma 4.7.** *For any $\varepsilon > 0$ and $N, L \in \mathbb{N}^+$, denote $M = N^2 L$ and $n = \lfloor \log_3(N+2) \rfloor$.*
*Assume $y_{m,k} \geq 0$ for $m = 0, 1, \cdots, M-1$ and $k = 0, 1, \cdots, Ln-1$ are samples with*

$$|y_{m,k} - y_{m,k-1}| \leq \varepsilon, \quad \text{for } m = 0, 1, \cdots, M-1 \quad \text{and} \quad k = 1, 2, \cdots, Ln-1.$$

*Then there exists $\phi \in \mathcal{NN}(\#\text{input} = 2; \text{ width} \leq 16N + 30; \text{ depth} \leq 5L + 7; \#\text{output} = 1)$*
*such that*

*(i) $|\phi(m, k) - y_{m,k}| \leq \varepsilon$ for $m = 0, 1, \cdots, M-1$ and $k = 0, 1, \cdots, Ln-1$;*

*(ii) $0 \leq \phi(x_1, x_2) \leq \max\{y_{m,k} : m = 0, 1, \cdots, M-1 \quad \text{and} \quad k = 0, 1, \cdots, Ln-1\}$ for any $x_1, x_2 \in \mathbb{R}$.*

*Proof.* Define

$$a_{m,k} := \lfloor y_{m,k}/\varepsilon \rfloor, \quad \text{for } m = 0, 1, \cdots, M-1 \quad \text{and} \quad k = 0, 1, \cdots, Ln-1.$$

We will construct a function implemented by a ReLU network to map the index $(m, k)$
to $a_{m,k}\varepsilon$ for $m = 0, 1, \cdots, M-1$ and $k = 0, 1, \cdots, Ln-1$.

Define $b_{m,0} := 0$ and $b_{m,k} := a_{m,k} - a_{m,k-1}$ for $m = 0, 1, \cdots, M-1$ and $k = 1, 2, \cdots, Ln-1$.
Since $|y_{m,k} - y_{m,k-1}| \leq \varepsilon$ for all $m$ and $k$, we have $b_{m,k} \in \{-1, 0, 1\}$. Hence, there exist
$c_{m,k} \in \{0, 1\}$ and $d_{m,k} \in \{0, 1\}$ such that $b_{m,k} = c_{m,k} - d_{m,k}$, which implies

$$a_{m,k} = a_{m,0} + \sum_{i=1}^{k}(a_{m,i} - a_{m,i-1}) = a_{m,0} + \sum_{i=1}^{k} b_{m,i} = a_{m,0} + \sum_{i=0}^{k} b_{m,i}$$

$$= a_{m,0} + \sum_{i=0}^{k} c_{m,i} - \sum_{i=0}^{k} d_{m,i},$$

for $m = 0, 1, \cdots, M-1$ and $k = 0, 1, \cdots, Ln-1$.

Consider the sample set

$$\left\{(m, a_{m,0}) : m = 0, 1, \cdots, M-1\right\} \bigcup \{(M, 0)\}.$$

Its size is $M + 1 = N \cdot ((NL - 1) + 1) + 1$, by Lemma 4.1 (set $N_1 = N$ and $N_2 = NL - 1$
therein), there exists

$$\psi_1 \in \mathcal{NN}(\text{widthvec} = [2N, 2(NL-1) + 1]) = \mathcal{NN}(\text{widthvec} = [2N, 2NL - 1])$$

such that
$$\psi_1(m) = a_{m,0}, \quad \text{for } m = 0, 1, \cdots, M-1.$$

By Lemma 4.6, there exist $\psi_2, \psi_3 \in \mathcal{NN}(\text{width} \leq 6N + 14; \text{ depth} \leq 5L + 4)$ such that

$$\psi_2(m, k) = \sum_{i=0}^{k} c_{m,i} \quad \text{and} \quad \psi_3(m, k) = \sum_{i=0}^{k} d_{m,i},$$

for $m = 0, 1, \cdots, M-1$ and $k = 0, 1, \cdots, Ln-1$. Hence, it holds that

$$a_{m,k} = a_{m,0} + \sum_{i=0}^{k} c_{m,i} - \sum_{i=0}^{k} d_{m,i} = \psi_1(m) + \psi_2(m, k) - \psi_3(m, k), \tag{4.11}$$

30

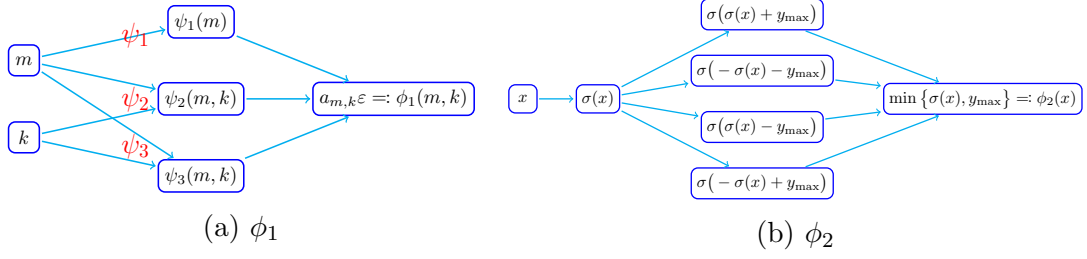(a) $\phi_1$                         (b) $\phi_2$

Figure 13: Illustrations of two sub-networks implementing the desired function $\phi = \phi_2 \circ \phi_1$ for $m = 0, 1, \cdots, M - 1$ and $k = 0, 1, \cdots, Ln - 1$, based on Equation (4.11) and the fact $\min\{x_1, x_2\} = \frac{x_1 + x_2 - |x_1 - x_2|}{2} = \frac{\sigma(x_1 + x_2) - \sigma(-x_1 - x_2) - \sigma(x_1 - x_2) - \sigma(-x_1 + x_2)}{2}$.

for $m = 0, 1, \cdots, M - 1$ and $k = 0, 1, \cdots, Ln - 1$.

Define

$$y_{\max} := \max\{y_{m,k} : m = 0, 1, \cdots, M - 1 \quad \text{and} \quad k = 0, 1, \cdots, Ln - 1\}.$$

Then the desired function can be implemented by two sub-networks shown in Figure 13.

By Lemma 4.2,

$$\psi_1 \in \mathcal{NN}(\#\text{input} = 1; \text{ widthvec} = [2N, 2NL - 1]; \#\text{output} = 1)$$
$$\subseteq \mathcal{NN}(\#\text{input} = 1; \text{ width} \leq 4N + 2; \text{ depth} \leq L + 1; \#\text{output} = 1).$$

Recall that $\psi_2, \psi_3 \in \mathcal{NN}(\text{width} \leq 6N + 14; \text{ depth} \leq 5L + 4)$. Thus, $\phi_1 \in \mathcal{NN}(\text{width} \leq (4N + 2) + 2(6N + 14) = 16N + 30; \text{ depth} \leq (5L + 4) + 1 = 5L + 5)$ as shown in Figure 13. And it is clear that $\phi_2 \in \mathcal{NN}(\text{width} \leq 4; \text{ depth} \leq 2)$, implying $\phi = \phi_2 \circ \phi_1 \in \mathcal{NN}(\text{width} \leq 16N + 30; \text{ depth} \leq (5L + 5) + 2 = 5L + 7)$.

Clearly, $0 \leq \phi(x_1, x_2) \leq y_{\max}$ for any $x_1, x_2 \in \mathbb{R}$, since $\phi(x_1, x_2) = \phi_2 \circ \phi_1(x_1, x_2) = \max\{\sigma(\phi_1(x_1, x_2)), y_{\max}\}$.

Note that $0 \leq a_{m,k}\varepsilon = \lfloor y_{m,k}/\varepsilon \rfloor \varepsilon \leq y_{\max}$. Then we have $\phi(m, k) = \phi_2 \circ \phi_1(m, k) = \phi_2(a_{m,k}\varepsilon) = \max\{\sigma(a_{m,k}\varepsilon), y_{\max}\} = a_{m,k}\varepsilon$. Therefore,

$$|\phi(m, k) - y_{m,k}| = |a_{m,k}\varepsilon - y_{m,k}| = |\lfloor y_{m,k}/\varepsilon \rfloor \varepsilon - y_{m,k}| \leq \varepsilon,$$

for $m = 0, 1, \cdots, M - 1$ and $k = 0, 1, \cdots, Ln - 1$. Hence, we finish the proof. $\qquad\square$

Finally, we apply Lemma 4.7 to prove Proposition 3.2.

*Proof of Proposition 3.2.* Denote $M = N^2 L$, $n = \lfloor \log_3(N + 2) \rfloor$, and $\widehat{L} = Ln$. We may assume $J = MLn = M\widehat{L}$ since we can set $y_{J-1} = y_J = y_{J+1} = \cdots = y_{M\widehat{L}-1}$ if $J < M\widehat{L}$.

Consider the sample set

$$\big\{(m\widehat{L}, m) : m = 0, 1, \cdots, M\big\} \bigcup \big\{(m\widehat{L} + \widehat{L} - 1, m) : m = 0, 1, \cdots, M - 1\big\}.$$

Its size is $2M + 1 = N \cdot \big((2NL - 1) + 1\big) + 1$. By Lemma 4.1 (set $N_1 = N$ and $N_2 = NL - 1$ therein), there exist

$$\phi_1 \in \mathcal{NN}(\text{widthvec} = [2N, 2(2NL - 1) + 1]) = \mathcal{NN}(\text{widthvec} = [2N, 4NL - 1])$$

such that

31

- $\phi_1(M\widehat{L}) = M$ and $\phi_1(m\widehat{L}) = \phi_1(m\widehat{L} + \widehat{L} - 1) = m$ for $m = 0, 1, \cdots, M - 1$.

- $\phi_1$ is linear on each interval $[m\widehat{L}, m\widehat{L} + \widehat{L} - 1]$ for $m = 0, 1, \cdots, M - 1$.

It follows that

$$\phi_1(j) = m, \quad \text{and} \quad j - \widehat{L}\phi_1(j) = k, \quad \text{where } j = m\widehat{L} + k, \tag{4.12}$$

for $m = 0, 1, \cdots, M - 1$ and $k = 0, 1, \cdots, \widehat{L} - 1$.

Note that any number $j$ in $\{0, 1, \cdots, J - 1\}$ can be uniquely indexed as $j = m\widehat{L} + k$ for $m = 0, 1, \cdots, M - 1$ and $k = 0, 1, \cdots, \widehat{L} - 1$. So we can denote $y_j = y_{m\widehat{L}+k}$ as $y_{m,k}$. Then by Lemma 4.7, there exists $\phi_2 \in \mathcal{NN}(\text{width} \le 16N + 30; \text{ depth} \le 5L + 7)$ such that

$$|\phi_2(m, k) - y_{m,k}| \le \varepsilon, \quad \text{for } m = 0, 1, \cdots, M - 1 \quad \text{and} \quad k = 0, 1, \cdots, \widehat{L} - 1, \tag{4.13}$$

and

$$0 \le \phi_2(x_1, x_2) \le y_{\max}, \quad \text{for any } x_1, x_2 \in \mathbb{R}, \tag{4.14}$$

where $y_{\max} := \max\{y_{m,k} : m = 0, 1, \cdots, M - 1 \text{ and } k = 0, 1, \cdots, \widehat{L} - 1\} = \max\{y_j : j = 0, 1, \cdots, M\widehat{L} - 1\}$.
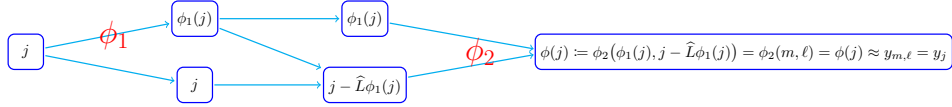


Figure 14: An illustration of the ReLU network implementing the desired function $\phi$ based Equation (4.12). The index $j \in \{0, 1, \cdots, M\widehat{L} - 1\}$ is unique represented by $j = mL + k$ for $m = 0, 1, \cdots, M - 1$ and $k = 0, 1, \cdots, \widehat{L} - 1$.

By Lemma 4.2,

$$\phi_1 \in \mathcal{NN}(\#\text{input} = 1; \text{ widthvec} = [2N, 4NL - 1]; \#\text{output} = 1)$$
$$\subseteq \mathcal{NN}(\#\text{input} = 1; \text{ width} \le 8N + 2; \text{depth} \le L + 1; \#\text{output} = 1).$$

Recall that $\phi_2 \in \mathcal{NN}(\text{width} \le 16N + 30; \text{ depth} \le 5L + 7)$. So $\phi \in \mathcal{NN}(\text{width} \le 16N + 30; \text{ depth} \le (L + 1) + 2 + (5L + 7) = 6L + 10)$ as shown in Figure 14.

Equation (4.14) implies

$$0 \le \phi(x) \le y_{\max}, \quad \text{for any } x \in \mathbb{R},$$

since $\phi$ is given by $\phi(x) = \phi_2\big(\phi_1(x), x - L\phi_1(x)\big)$.

Represent $j \in \{0, 1, \cdots, M\widehat{L} - 1\}$ via $j = mL + k$ for $m = 0, 1, \cdots, M - 1$ and $k = 0, 1, \cdots, \widehat{L} - 1$. Then, by Equation (4.13), we have

$$|\phi(j) - y_j| = |\phi_2\big(\phi_1(j), j - L\phi_1(j)\big) - y_j| = |\phi_2(m, k) - y_{m,k}| \le \varepsilon,$$

for any $j \in \{0, 1, \cdots, M\widehat{L} - 1\} = \{0, 1, \cdots, J - 1\}$. So we finish the proof. $\square$

We would like to remark that the key idea in the proof of Proposition 3.2 is the bit extraction technique in Lemma 4.5, which allows us to store $Ln$ bits in a binary number $\text{bin}\,0.\theta_1\theta_2\cdots\theta_{Ln}$ and extract each bit $\theta_i$. The extraction operator can be efficiently carried out via a deep ReLU neural network demonstrating the power of depth.

# 5 Conclusion and future work

This paper aims at a quantitative and optimal approximation rate for ReLU networks in terms of the width and depth to approximate continuous functions. It is shown by construction that ReLU networks with width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ can approximate an arbitrary continuous function on $[0,1]^d$ with an approximation rate $\mathcal{O}\big(\omega_f\big((N^2L^2\ln N)^{-1/d}\big)\big)$. By connecting the approximation property to VC-dimension, we prove that such a rate is optimal for Hölder continuous functions on $[0,1]^d$ in terms of the width and depth separately, and hence this rate is also optimal for the whole continuous function class. We also extend our analysis to general continuous functions on any bounded set in $\mathbb{R}^d$. We would like to remark that our analysis was based on the fully connected feed-forward neural networks and the ReLU activation function. It would be very interesting to extend our conclusions to neural networks with other types of architectures (e.g., convolutional neural networks) and activation functions (e.g., tanh and sigmoid functions).

# References

[1] M. ANTHONY AND P. L. BARTLETT, *Neural Network Learning: Theoretical Foundations*, Cambridge University Press, New York, NY, USA, 1st ed., 2009.

[2] A. R. BARRON, *Universal approximation bounds for superpositions of a sigmoidal function*, IEEE Transactions on Information Theory, 39 (1993), pp. 930–945.

[3] P. BARTLETT, V. MAIOROV, AND R. MEIR, *Almost linear VC dimension bounds for piecewise polynomial networks*, Neural Computation, 10 (1998), pp. 2159–2173.

[4] M. BIANCHINI AND F. SCARSELLI, *On the complexity of neural network classifiers: A comparison between shallow and deep architectures*, IEEE Transactions on Neural Networks and Learning Systems, 25 (2014), pp. 1553–1565.

[5] G. CYBENKO, *Approximation by superpositions of a sigmoidal function*, MCSS, 2 (1989), pp. 303–314.

[6] W. E AND Q. WANG, *Exponential convergence of the deep neural network approximation for analytic functions*, CoRR, abs/1807.00297 (2018).

33

[7] W. E AND S. WOJTOWYTSCH, *On the banach spaces associated with multi-layer relu networks: Function representation, approximation theory and gradient descent dynamics*, arXiv:2007.15623, (2020).

[8] N. HARVEY, C. LIAW, AND A. MEHRABIAN, *Nearly-tight VC-dimension bounds for piecewise linear neural networks*, in Proceedings of the 2017 Conference on Learning Theory, S. Kale and O. Shamir, eds., vol. 65 of Proceedings of Machine Learning Research, Amsterdam, Netherlands, 07–10 Jul 2017, PMLR, pp. 1064–1068.

[9] K. HORNIK, M. STINCHCOMBE, AND H. WHITE, *Multilayer feedforward networks are universal approximators*, Neural Networks, 2 (1989), pp. 359 – 366.

[10] K. KAWAGUCHI, *Deep learning without poor local minima*, in Advances in Neural Information Processing Systems 29, D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, eds., Curran Associates, Inc., 2016, pp. 586–594.

[11] K. KAWAGUCHI AND Y. BENGIO, *Depth with nonlinearity creates no bad local minima in resnets*, (2018).

[12] M. J. KEARNS AND R. E. SCHAPIRE, *Efficient distribution-free learning of probabilistic concepts*, J. Comput. Syst. Sci., 48 (1994), pp. 464–497.

[13] Q. LI, T. LIN, AND Z. SHEN, *Deep learning via dynamical systems: An approximation perspective*, arXiv e-prints, (2019).

[14] S. LIANG AND R. SRIKANT, *Why deep neural networks?*, CoRR, abs/1610.04161 (2016).

[15] J. LU, Z. SHEN, H. YANG, AND S. ZHANG, *Deep network approximation for smooth functions*, arXiv e-prints, (2020), p. arXiv:2001.03040.

[16] Z. LU, H. PU, F. WANG, Z. HU, AND L. WANG, *The expressive power of neural networks: A view from the width*, in Advances in Neural Information Processing Systems 30, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, eds., Curran Associates, Inc., 2017, pp. 6231–6239.

[17] H. MONTANELLI, H. YANG, AND Q. DU, *Deep ReLU networks overcome the curse of dimensionality for bandlimited functions*, Journal of Computational Mathematics, (to appear).

[18] G. F. MONTUFAR, R. PASCANU, K. CHO, AND Y. BENGIO, *On the number of linear regions of deep neural networks*, in Advances in Neural Information Processing Systems 27, Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, eds., Curran Associates, Inc., 2014, pp. 2924–2932.

[19] Q. N. NGUYEN AND M. HEIN, *The loss surface of deep and wide neural networks*, CoRR, abs/1704.08045 (2017).

[20] P. Petersen and F. Voigtlaender, *Optimal approximation of piecewise smooth functions using deep ReLU neural networks*, Neural Networks, 108 (2018), pp. 296 – 330.

[21] A. Sakurai, *Tight bounds for the VC-dimension of piecewise polynomial networks*, in Advances in Neural Information Processing Systems, Neural information processing systems foundation, 1999, pp. 323–329.

[22] Z. Shen, H. Yang, and S. Zhang, *Nonlinear approximation via compositions*, Neural Networks, 119 (2019), pp. 74 – 84.

[23] ——, *Deep network approximation characterized by number of neurons*, Communications in Computational Physics, 28 (2020), pp. 1768–1811.

[24] Z. Shen, H. Yang, and S. Zhang, *Deep network with approximation error being reciprocal of width to power of square root of depth*, arXiv e-prints, (2020).

[25] ——, *Neural network approximation: Three hidden layers are enough*, arXiv e-prints, (2020).

[26] J. W. Siegel and J. Xu, *Approximation rates for neural networks with general activation functions*, Neural Networks, 128 (2020), pp. 313–321.

[27] D. Yarotsky, *Error bounds for approximations with deep ReLU networks*, Neural Networks, 94 (2017), pp. 103 – 114.

[28] D. Yarotsky, *Optimal approximation of continuous functions by very deep ReLU networks*, in Proceedings of the 31st Conference On Learning Theory, S. Bubeck, V. Perchet, and P. Rigollet, eds., vol. 75 of Proceedings of Machine Learning Research, PMLR, 06–09 Jul 2018, pp. 639–649.

[29] S. Zhang, *Deep neural network approximation via function compositions*, PhD Thesis, National University of Singapore, (2020). URL: https://scholarbank.nus.edu.sg/handle/10635/186064.