

# ReLU Network Approximation in Terms of Intrinsic Parameters

Zuowei Shen\*      Haizhao Yang†      Shijun Zhang‡

## Abstract

This paper studies the approximation error of ReLU networks in terms of the number of intrinsic parameters (i.e., those depending on the target function  $f$ ). First, we prove by construction that, for any Lipschitz continuous function  $f$  on  $[0, 1]^d$  with a Lipschitz constant  $\lambda > 0$ , a ReLU network with  $n + 2$  intrinsic parameters can approximate  $f$  with an exponentially small error  $5\lambda\sqrt{d}2^{-n}$  measured in the  $L^p$ -norm for  $p \in [1, \infty)$ . More generally for an arbitrary continuous function  $f$  on  $[0, 1]^d$  with a modulus of continuity  $\omega_f(\cdot)$ , the approximation error is  $\omega_f(\sqrt{d}2^{-n}) + 2^{-n+2}\omega_f(\sqrt{d})$ . Next, we extend these two results from the  $L^p$ -norm to the  $L^\infty$ -norm at a price of  $3^d n + 2$  intrinsic parameters. Finally, by using a high-precision binary representation and the bit extraction technique via a fixed ReLU network independent of the target function, we design, theoretically, a ReLU network with only three intrinsic parameters to approximate Hölder continuous functions with an arbitrarily small error.

**Key words:** ReLU Neural Networks; Intrinsic Parameters; Approximation Error; Exponential Convergence; Transfer Learning.

## 1 Introduction

Deep learning has been a powerful tool in science and engineering. As the workhorses of deep learning, deep neural networks are treated as an important regression tool in many successful applications. Understanding the approximation capacity of deep neural networks has become a key question for revealing the power of deep learning. One of the fundamental problems is the characterization of the approximation error of deep neural networks in terms of the network size measured in the width, the depth, the number of neurons, or the number of parameters.

It was shown in [38, 42, 45, 48] that the approximation error  $\mathcal{O}(n^{-2/d})$  is (nearly) optimal for ReLU networks with  $\mathcal{O}(n)$  parameters to approximate Lipschitz continuous functions on  $[0, 1]^d$ . To gain better approximation errors, existing results either consider smaller target function spaces (e.g., [3, 11, 15, 25, 29, 44, 47]) or introduce new activation

---

\*Department of Mathematics, National University of Singapore ([matzuows@nus.edu.sg](mailto:matzuows@nus.edu.sg)).

†Department of Mathematics, Purdue University ([haizhao@purdue.edu](mailto:haizhao@purdue.edu)).

‡Department of Mathematics, National University of Singapore ([zhangshijun@u.nus.edu](mailto:zhangshijun@u.nus.edu)).

functions (e.g., [39, 40, 41, 46]). This paper proposes a new perspective to study the approximation error in terms of the number of parameters depending on the target function, which are called the **intrinsic parameters**, excluding those independent of the target function.

## 1.1 Main results

Let  $C([0, 1]^d)$  denote the space of continuous functions defined on  $[0, 1]^d$ . For simplicity, let  $\mathcal{H}_W(d_1, d_2)$  denote the function space consisting of all ReLU networks with  $W$  parameters mapping from  $\mathbb{R}^{d_1}$  to  $\mathbb{R}^{d_2}$ , i.e.,

$$\mathcal{H}_W(d_1, d_2) := \left\{ g : g : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_2} \text{ is realized by a ReLU network with } W \text{ parameters} \right\}.$$

Let  $\mathcal{H}(d_1, d_2) := \bigcup_{W=1}^{\infty} \mathcal{H}_W(d_1, d_2)$ .

For any  $f \in C([0, 1]^d)$ , our goal is to construct two  $f$ -independent functions  $\phi_1 \in \mathcal{H}(d, 1)$  and  $\phi_2 \in \mathcal{H}(n, 1)$ , and use  $s \cdot (\phi_2 \circ \phi_f \circ \phi_1) + b$  to approximate  $f$ , where  $s \in [0, \infty)$ ,  $b \in \mathbb{R}$ , and  $\phi_f \in \mathcal{H}_n(1, n)$  are learned from  $f$ . Under these settings, an approximation error  $\omega_f(\sqrt{d}2^{-n}) + 2^{-n+2}\omega_f(\sqrt{d})$  is attained as shown in the theorem below, where the modulus of continuity of a continuous function  $f \in C([0, 1]^d)$  is defined as

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

**Theorem 1.1.** *Given any  $n \in \mathbb{N}^+$  and  $p \in [1, \infty)$ , there exist  $\phi_1 \in \mathcal{H}_{2^{dn+4}}(d, 1)$  and  $\phi_2 \in \mathcal{H}_{2^{dn+5n}}(n, 1)$  such that: For any  $f \in C([0, 1]^d)$ , there exists a linear map  $\mathcal{L} : \mathbb{R} \rightarrow \mathbb{R}^n$  satisfying*

$$\left\| s \cdot (\phi_2 \circ \mathcal{L} \circ \phi_1) + b - f \right\|_{L^p([0, 1]^d)} \leq \omega_f(\sqrt{d}2^{-n}) + 2^{-n+2}\omega_f(\sqrt{d}),$$

where  $s = 2\omega_f(\sqrt{d})$ ,  $b = f(\mathbf{0}) - \omega_f(\sqrt{d})$ , and  $\mathcal{L}$  is a linear map given by  $\mathcal{L}(t) = (a_1t, a_2t, \dots, a_nt)$  with  $a_1, a_2, \dots, a_n \in [0, \frac{1}{3})$  determined by  $f$  and  $n$ .

In Theorem 1.1,  $s$  is a scale factor,  $b$  is the bias for a vertical shift, and  $a_1, a_2, \dots, a_n \in [0, \frac{1}{3})$  are the key intrinsic parameters storing most of information of  $f$ . Clearly,  $s \cdot (\phi_2 \circ \mathcal{L} \circ \phi_1) + b$  can be implemented by a ReLU network with  $n + 2$  intrinsic parameters. We call  $\phi_1$  and  $\phi_2$  **inner-function** and **outer-function**, respectively. They are independent of the target function  $f$  and can be implemented by ReLU networks.

Note that the approximation error in Theorem 1.1 is characterized by the  $L^p$ -norm for  $p \in [1, \infty)$ . In fact, we can extend such a result to a similar one measured in the  $L^\infty$ -norm.

**Theorem 1.2.** *Given any  $n \in \mathbb{N}^+$ , there exist  $\phi_1 \in \mathcal{H}_{3^{d2^{dn+5}}}(d, 3^d)$  and  $\phi_2 \in \mathcal{H}_{3^{d2^{dn+8n}}}(3^dn, 1)$  such that: For any  $f \in C([0, 1]^d)$ , there exists a linear map  $\mathcal{L} : \mathbb{R}^{3^d} \rightarrow \mathbb{R}^{3^dn}$  satisfying*

$$\left\| s \cdot (\phi_2 \circ \mathcal{L} \circ \phi_1) + b - f \right\|_{L^\infty([0, 1]^d)} \leq \omega_f(\sqrt{d}2^{-n}) + 2^{-n+2}\omega_f(\sqrt{d}),$$

where  $s = 2\omega_f(\sqrt{d})$ ,  $b = f(\mathbf{0}) - \omega_f(\sqrt{d})$ , and  $\mathcal{L}$  is given by

$$\mathcal{L}(x_1, \dots, x_{3^d}) = \left( \mathcal{L}_0(x_1), \dots, \mathcal{L}_0(x_{3^d}) \right) \quad \text{for any } \mathbf{x} = (x_1, \dots, x_{3^d}) \in \mathbb{R}^{3^d},$$

where  $\mathcal{L}_0 : \mathbb{R} \rightarrow \mathbb{R}^n$  is a linear map given by  $\mathcal{L}_0(t) = (a_1t, a_2t, \dots, a_nt)$  with  $a_1, a_2, \dots, a_n \in [0, \frac{1}{3})$  determined by  $f$  and  $n$ .

Simplifying the implicit approximation error in Theorem 1.1 (or 1.2) to make it explicitly depending on  $n$  is challenging in general, since the modulus of continuity  $\omega_f(\cdot)$  may be complicated. However, if  $f$  is a Hölder continuous function on  $[0, 1]^d$  of order  $\alpha \in (0, 1]$  with a Hölder constant  $\lambda > 0$ . That is,  $f$  satisfies

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq \lambda \|\mathbf{x} - \mathbf{y}\|_2^\alpha \quad \text{for any } \mathbf{x}, \mathbf{y} \in [0, 1]^d,$$

implying  $\omega_f(r) \leq \lambda r^\alpha$  for any  $r \geq 0$ . This means we can get an exponentially small approximation error  $5\lambda d^{\alpha/2} 2^{-\alpha n}$ . In particular, in the special case of  $\alpha = 1$ , i.e.,  $f$  is a Lipschitz continuous function with a Lipschitz constant  $\lambda > 0$ , then the approximation error is simplified to  $5\lambda\sqrt{d}2^{-n}$ .

Though the linear map  $\mathcal{L}$  in Theorem 1.2 is essentially determined by  $n$  key parameters  $a_1, a_2, \dots, a_n$ , these  $n$  key parameters are repeated  $3^d$  times in the final network architecture as shown in Figure 3. Therefore,  $s \cdot (\phi_2 \circ \mathcal{L} \circ \phi_1) + b$  can be implemented by a ReLU network with  $3^d n + 2$  intrinsic parameters. Remark that we can reduce the number of intrinsic parameters to  $n + 2$  via using a fixed ReLU network to copy  $n$  key parameters  $3^d$  times. The idea is similar to that of Theorem 1.3 introduced later.

Furthermore, the number of intrinsic parameters can be reduced to three in the case of Hölder continuous functions. In fact, three intrinsic parameters are enough to achieve an arbitrary pre-specified error if sufficiently high precision is provided, as shown in the theorem below.

**Theorem 1.3.** *Given any  $\varepsilon > 0$ ,  $\alpha \in (0, 1]$ , and  $\lambda > 0$ , there exists  $\phi \in \mathcal{H}(d + 1, 1)$  such that: For any Hölder continuous function  $f$  on  $[0, 1]^d$  of order  $\alpha \in (0, 1]$  with a Hölder constant  $\lambda > 0$ , there exist three parameters  $s \in [0, \infty)$ ,  $v \in [0, 1)$ , and  $b \in \mathbb{R}$  satisfying*

$$|s\phi(\mathbf{x}, v) + b - f(\mathbf{x})| \leq \varepsilon \quad \text{for any } \mathbf{x} \in [0, 1]^d.$$

In Theorem 1.3,  $s$  is a scale factor,  $b$  is the bias for a vertical shift, and  $v$  is the key intrinsic parameter storing sufficient information of the target function  $f$ . Clearly,  $s$ ,  $b$ , and  $v$  are learned from  $f$ , while  $\phi$  is independent of  $f$ . Let  $\phi_2 = \phi$ ,  $\phi_1$  be the identity map on  $\mathbb{R}^d$ , and  $\mathcal{L}_v : \mathbb{R}^d \rightarrow \mathbb{R}^{d+1}$  be an affine linear transform mapping  $\mathbf{x}$  to  $(\mathbf{x}, v)$ . Then  $s\phi(\mathbf{x}, v) + b$  can also be represented as  $s\phi_2 \circ \mathcal{L}_v \circ \phi_1(\mathbf{x}) + b$ .

Remark that Theorem 1.3 is just a theoretical result since the key intrinsic parameter “ $v$ ” requires extremely high precision, which is necessary for storing the values of  $f$  at sufficiently many points within a sufficiently small error. Via the idea of the binary representation, we can extract the values of  $f$  stored in “ $v$ ” via an  $f$ -independent ReLU network (as a sub-network of the final network realizing  $\phi$  in Theorem 1.3). In fact, there is a balance between the precision requirement and the number of intrinsic parameters. For example, if we store the values of  $f$  in two intrinsic parameters (not one), then the precision requirement is greatly lessened.

## 1.2 Contributions and further interpretation

Our key contributions can be summarized as follows.

- (i) First, we prove by construction in Theorem 1.1 that a ReLU network with  $n + 2$  intrinsic parameters can approximate a continuous function  $f$  on  $[0, 1]^d$  with an

error  $\omega_f(\sqrt{d}2^{-n}) + 2^{-n+2}\omega_f(\sqrt{d})$  measured in the  $L^p$ -norm for  $p \in [1, \infty)$ . In the case of Hölder continuous functions, the approximation error is simplified to  $5\lambda d^{\alpha/2}2^{-\alpha n}$ , where  $\alpha \in (0, 1]$  and  $\lambda > 0$  are the Hölder order and constant, respectively.

- (ii) Next, we generalize the approximation error in Theorem 1.1 from the  $L^p$ -norm for  $p \in [1, \infty)$  to the  $L^\infty$ -norm, as shown in Theorem 1.2. Such a generalization is at a price of more intrinsic parameters. To be precise, the final network in Theorem 1.2 has  $3^d n + 2$  intrinsic parameters, compared to  $n + 2$  ones in Theorem 1.1.
- (iii) Finally, we show in Theorem 1.3 that the number of intrinsic parameters in Theorems 1.1 and 1.2 can be further reduced to three in the case of Hölder continuous functions. To be precise, ReLU networks with three intrinsic parameters can achieve an arbitrary error for approximating Hölder continuous function on  $[0, 1]^d$ . In this scenario, extremely high precision is required as we shall see later.

### Approximation with inner-function and outer-function

The composition architecture  $\phi_2 \circ \mathcal{L} \circ \phi_1$  is the key part of the final ReLU-network-realized function  $s \cdot (\phi_2 \circ \mathcal{L} \circ \phi_1) + b$  in Theorem 1.1 (or 1.2). Given a target function  $f$ , the composition architecture  $\phi_2 \circ \mathcal{L} \circ \phi_1$  can be generalized to  $\phi_{\text{out}} \circ \phi_\theta \circ \phi_{\text{in}}$ , where  $\phi_\theta$  an  $f$ -dependent function parameterized by  $\theta = \theta(f) \in \mathbb{R}^n$ , and  $\phi_{\text{in}}$  and  $\phi_{\text{out}}$  are two  $f$ -independent functions, called the inner-function and the outer-function, respectively.

Let  $\Gamma_n(d_1, d_2)$  denote a general space consisting of vector-valued functions mapping from  $\mathbb{R}^{d_1}$  to  $\mathbb{R}^{d_2}$  and parameterized with  $n$  parameters. Each element of  $\Gamma_n(d_1, d_2)$  can be denoted by  $\phi_\theta : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_2}$ , parameterized with  $\theta \in \mathbb{R}^n$ . Existing literature uses  $\phi_\theta \in \Gamma_n(d_1, d_2)$  with proper  $d_1$  and  $d_2$  to directly approximate the target function  $f$  in a given function space. This paper proposes a new perspective to study the function approximation for  $f : [0, 1]^d \rightarrow \mathbb{R}$ . To be precise, we design an inner-function  $\phi_{\text{in}} : \mathbb{R}^d \rightarrow \mathbb{R}^{d_1}$  and an outer-function  $\phi_{\text{out}} : \mathbb{R}^{d_2} \rightarrow \mathbb{R}$ , both of which are independent of  $f$ , and use  $s \cdot (\phi_{\text{out}} \circ \phi_\theta \circ \phi_{\text{in}}) + b$  with  $\phi_\theta \in \Gamma_n(d_1, d_2)$  to approximate the target function  $f$ . Here,  $s$  is a scale factor,  $b$  is the bias for a vertical shift, and  $\theta \in \mathbb{R}^n$  is the key parameter vector learned from  $f$ .

In particular, let  $\Gamma_n$  denote the function of all functions realized by a ReLU network with  $n$  parameters and a pre-specified architecture. It is proved in [38, 45] that, when using elements in  $\Gamma_n$  to approximate Hölder continuous functions on  $[0, 1]^d$ , the (nearly) optimal approximation error is  $\mathcal{O}(\lambda n^{-2\alpha/d})$ , where  $\alpha \in (0, 1]$  and  $\lambda > 0$  are the Hölder order and constant, respectively. Clearly, an error like  $\mathcal{O}(\lambda n^{-2\alpha/d})$  suffers from the curse of dimensionality. However, with inner-function  $\phi_{\text{in}}$  and outer-function  $\phi_{\text{out}}$  in hand, we can use  $s \cdot (\phi_{\text{out}} \circ \phi_\theta \circ \phi_{\text{in}}) + b$  with  $\phi_\theta \in \Gamma_n$  to approximate a Hölder continuous function on  $[0, 1]^d$  with an exponentially small error  $5\lambda d^{\alpha/2}2^{-\alpha n}$ , where  $\alpha \in (0, 1]$  and  $\lambda > 0$  are the Hölder order and constant, respectively. This means that the approximation error can be greatly improved by pre-designing two (vector-valued) functions  $\phi_{\text{in}}$  and  $\phi_{\text{out}}$ , which are independent of  $f$  and can be implemented by ReLU networks.

## Connection with transfer learning

Transfer learning dates back to 1970s [9, 10]. It is a research direction in machine learning that applies knowledge gained in one problem to solve a different but related problem. Typically in deep learning, transfer learning uses a pre-trained neural network obtained for one task as an initial guess of the neural network for another task to achieve a short training time. Our theory in this paper could provide insights into the success of transfer learning using neural networks, though the setting of our theory is different from realistic transfer learning. In our theory,  $\phi_{\text{in}}$  and  $\phi_{\text{out}}$  are universally useful for all learning tasks for continuous functions, which can be understood as the part of networks that can be transferred to different tasks. Suppose  $f_1$  and  $f_2$  are the target functions for two different but related tasks. If  $f_1$  has been learned via an architecture  $s_1 \cdot (\phi_{\text{out}} \circ \phi_{\theta_1} \circ \phi_{\text{in}}) + b_1$ , then we can “transfer” the knowledge ( $\phi_{\text{out}}$  and  $\phi_{\text{in}}$ ) to another task. This means that, by only learning  $s_2, b_2, \theta_2$  from  $f_2$ , we can use  $s_2 \cdot (\phi_{\text{out}} \circ \phi_{\theta_2} \circ \phi_{\text{in}}) + b_2$  to approximate  $f_2$  well. Therefore, the total number of parameters that need to be learned again is not large. Our theory may provide a certain theoretical understanding in the spirit of transfer learning from a network approximation perspective. To gain a deeper understanding, one can refer to [6, 22, 26, 27, 32, 34, 35]. It would be interesting to test the proposed network architecture in the context of transfer learning in the future.

## Error analysis of deep learning

In supervised learning, an unknown target function  $f$  defined on a domain  $\mathcal{X}$  is learned through its finitely many samples  $\{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n$ . For simplicity, denote  $\phi(\mathbf{x}; \boldsymbol{\theta})$  as a network-generated function with  $\boldsymbol{\theta} \in \Theta$  as the set of parameters, where  $\Theta$  is the parameter domain typically taken as  $[-M, M]^W$  for two pre-specified constants  $M > 0$  and  $W \in \mathbb{N}^+$ . If  $\phi(\mathbf{x}; \boldsymbol{\theta})$  is used to infer  $f(\mathbf{x})$  for unseen data samples  $\mathbf{x}$ , then we need to identify the empirical risk minimizer  $\boldsymbol{\theta}_S$ , which is given by

$$\boldsymbol{\theta}_S \in \arg \min_{\boldsymbol{\theta} \in \Theta} R_S(\boldsymbol{\theta}), \quad \text{where } R_S(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n \ell(\phi(\mathbf{x}_i; \boldsymbol{\theta}), f(\mathbf{x}_i)) \quad (1.1)$$

with a loss function typically taken as  $\ell(y, y') = \frac{1}{2}|y - y'|^2$ .

In fact, the best network-generated function to infer  $f(\mathbf{x})$  is  $\phi(\mathbf{x}; \boldsymbol{\theta}_D)$ , but not  $\phi(\mathbf{x}; \boldsymbol{\theta}_S)$ , where  $\boldsymbol{\theta}_D$  is the expected risk minimizer given by

$$\boldsymbol{\theta}_D \in \arg \min_{\boldsymbol{\theta} \in \Theta} R_D(\boldsymbol{\theta}), \quad \text{where } R_D(\boldsymbol{\theta}) := \mathbb{E}_{\mathbf{x} \sim U(\mathcal{X})} [\ell(\phi(\mathbf{x}; \boldsymbol{\theta}), f(\mathbf{x}))],$$

where  $U$  is a unknown data distribution over  $\mathcal{X}$ . The best possible inference error is  $R_D(\boldsymbol{\theta}_D)$ . In real applications,  $U(\mathcal{X})$  is unknown and only finitely many samples from this distribution are available. Hence, the empirical risk  $R_S(\boldsymbol{\theta})$  is minimized, hoping to obtain  $\phi(\mathbf{x}; \boldsymbol{\theta}_S)$ , instead of minimizing the expected risk  $R_D(\boldsymbol{\theta})$  to obtain  $\phi(\mathbf{x}; \boldsymbol{\theta}_D)$ . In practice, a numerical optimization method to solve (1.1) may result in a numerical solution (denoted as  $\boldsymbol{\theta}_N$ ) that may not be a global minimizer  $\boldsymbol{\theta}_S$ . Therefore, the actually learned network-generated function to infer  $f(\mathbf{x})$  is  $\phi(\mathbf{x}; \boldsymbol{\theta}_N)$  and the corresponding

inference error is measured by  $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}})$ , which is bounded by

$$\begin{aligned}
 R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}}) &= \underbrace{[R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}})]}_{\text{GE}} + \underbrace{[R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}})]}_{\text{OE}} + \underbrace{[R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{D}})]}_{\leq 0 \text{ by (1.1)}} + \underbrace{[R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{D}}) - R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})]}_{\text{GE}} + \underbrace{R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})}_{\text{AE}} \\
 &\leq \underbrace{R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})}_{\text{Approximation error (AE)}} + \underbrace{[R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}})]}_{\text{Optimization error (OE)}} + \underbrace{[R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}})] + [R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{D}}) - R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})]}_{\text{Generalization error (GE)}}. \quad (1.2)
 \end{aligned}$$

The constructive approximation established in this paper and the literature provides an upper bound of  $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})$ . The second term of (1.2) is bounded by the optimization error of the numerical algorithm applied to solve the empirical risk minimization problem in (1.1). If the numerical algorithm can find a global minimizer, the second term is equal to zero. The theoretical guarantee of the convergence of an optimization algorithm to a global minimizer  $\boldsymbol{\theta}_{\mathcal{S}}$  and the characterization of the convergence belong to the optimization analysis of neural networks. The study of the bounds for the third and fourth terms of (1.2) is referred to as the generalization error analysis of neural networks.

Theorems 1.1, 1.2, and 1.3 provide upper bounds of  $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})$ . These bounds only depends on the number of intrinsic parameters of ReLU networks and the modulus of continuity  $\omega_f(\cdot)$ . Hence, these bounds are independent of the empirical risk minimization in (1.1) and the optimization algorithm used to compute the numerical solution of (1.1). In other words, Theorems 1.1, 1.2, and 1.3 quantify the approximation power of ReLU networks in terms of the nubmer of intrinsic parameters. Designing efficient optimization algorithms and analyzing the generalization erre for ReLU networks are two other separate future directions.

### 1.3 Related work

The expressiveness of deep neural networks has been studied extensively from many perspectives, e.g., in terms of combinatorics [30], topology [7], Vapnik-Chervonenkis (VC) dimension [5, 18, 36], fat-shattering dimension [1, 23], information theory [33], classical approximation theory [2, 3, 8, 12, 13, 14, 16, 17, 20, 24, 25, 28, 31, 37, 38, 43, 44, 45, 48, 49], etc. In the early works of approximation theory for neural networks, the universal approximation theorem [14, 19, 20] without approximation errors showed that, given any  $\varepsilon > 0$ , there exists a sufficiently large neural network approximating a target function in a certain function space within an error  $\varepsilon$ . For one-hidden-layer neural networks and sufficiently smooth functions, Barron [3, 4] showed an asymptotic approximation error  $\mathcal{O}(\frac{1}{\sqrt{N}})$  in the  $L^2$ -norm, leveraging an idea that is similar to Monte Carlo sampling for high-dimensional integrals.

Recently, it is proved in [38, 45, 48] that the (nearly) optimal approximation error would be  $\mathcal{O}(n^{-2/d})$  when using ReLU networks with  $n$  parameters to approximate functions in the unit ball of Lipschitz continuous function space. Clearly, such an error suffers from the curse of dimensionality. To bridge this gap, one could either consider smaller function spaces, e.g., smooth functions [25, 47] and band-limited functions [29], or introducing new networks, e.g., Floor-ReLU networks [40], Floor-Exponential-Step (FLES) networks [41], and (Sin, ReLU,  $2^x$ )-activated networks [21]. This paper proposes a new perspective to characterize the approximation error in terms of the number of intrinsic parameters. Such a method is inspired by an observation that most parameters of the



ReLU network approximating a target function are independent of the target function. Thus, most parameters can be assigned or computed in advance. As shown in Theorem 1.1, we can first design an inner-function  $\phi_1$  and an outer-function  $\phi_2$ , both of which can be implemented by ReLU networks. Then, for any continuous function  $f \in C([0, 1]^d)$ ,  $s \cdot (\phi_2 \circ \mathcal{L} \circ \phi_1) + b$  can approximate  $f$  with an error  $\omega_f(\sqrt{d}2^{-n}) + 2^{-n+2}\omega_f(\sqrt{d})$  by the following two steps: 1) determining  $s$  and  $b$ , 2) designing a linear map  $\mathcal{L}$  defined by  $\mathcal{L}(t) = (a_1t, \dots, a_nt)$ , where  $a_1, \dots, a_n$  are determined by the target function  $f$ . Therefore, we overcome the curse of dimensionality in the sense of the approximation error characterized by the number of intrinsic parameters when the variation of  $\omega_f(r)$  as  $r \rightarrow 0$  is moderate (e.g.,  $\omega_f(r) \lesssim r^\alpha$  for Hölder continuous functions).

**Organization:** The rest of this paper is organized as follows. In Section 2, we prove Theorems 1.1, 1.2, and 1.3 based on an auxiliary theorem, Theorem 2.1. Next, the auxiliary theorem is proved in Section 3 based on Proposition 3.1, which is proved later at the end of Section 3. Finally, Section 4 concludes this paper with a short discussion.

## 2 Constructive proof

In this section, we first list all notations used throughout this paper. Then, we prove Theorems 1.1, 1.2, and 1.3 based on an auxiliary theorem, Theorem 2.1, which will be proved in Section 3.

### 2.1 Notations

Firstly, let us summarize the main notations of this paper as follows.

- Let  $\mathbb{R}$ ,  $\mathbb{Q}$ , and  $\mathbb{Z}$  denote the set of real numbers, rational numbers, and integers, respectively.
- Let  $\mathbb{N}$  and  $\mathbb{N}^+$  denote the set of natural numbers and positive natural numbers, respectively. That is,  $\mathbb{N}^+ = \{1, 2, 3, \dots\}$  and  $\mathbb{N} = \mathbb{N}^+ \cup \{0\}$ .
- Vectors and matrices are denoted in a bold font. Standard vectorization is adopted in the matrix and vector computation. For example, adding a scalar and a vector means adding the scalar to each entry of the vector.
- For  $\theta \in [0, 1)$ , 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  $\text{bin}0.\theta_1\theta_2\cdots\theta_L$  to denote the  $L$ -term binary representation of  $\theta$ , i.e.,  $\text{bin}0.\theta_1\theta_2\cdots\theta_L := \sum_{\ell=1}^L \theta_\ell 2^{-\ell}$ .
- For any  $p \in [1, \infty)$ , the  $p$ -norm of a vector  $\mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$  is defined by

$$\|\mathbf{x}\|_p := \left(|x_1|^p + |x_2|^p + \dots + |x_d|^p\right)^{1/p}.$$

- 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  $L$ .

- Similar to “min” and “max”, let  $\text{mid}(x_1, x_2, x_3)$  be the middle value of three inputs  $x_1$ ,  $x_2$ , and  $x_3$ . For example,  $\text{mid}(2, 1, 3) = 2$  and  $\text{mid}(3, 2, 3) = 3$ .
- 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

$$\Omega([0, 1]^d, K, \delta) := \bigcup_{i=1}^d \left\{ \mathbf{x} = (x_1, x_2, \dots, x_d) : x_i \in \bigcup_{k=1}^{K-1} \left( \frac{k}{K} - \delta, \frac{k}{K} \right) \right\}. \quad (2.1)$$

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

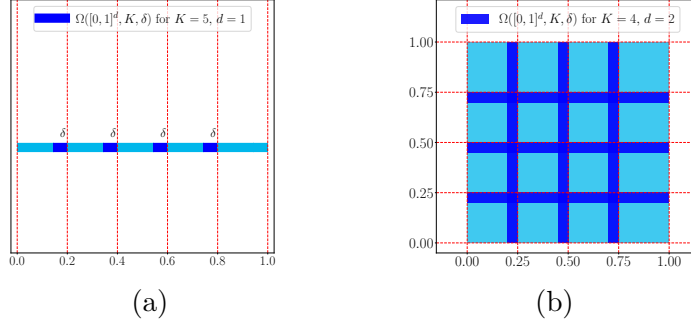


Figure 1: Two examples of trifling regions. (a)  $K = 5, d = 1$ . (b)  $K = 4, d = 2$ .

- Given a univariate activation function  $\sigma$ , let us introduce the architecture of a  $\sigma$ -activated network, i.e., a network with each hidden neuron activated by  $\sigma$ . To be precise, a  $\sigma$ -activated network with a vector input  $\mathbf{x} \in \mathbb{R}^d$ , an output  $\phi(\mathbf{x}) \in \mathbb{R}$ , and  $L \in \mathbb{N}^+$  hidden layers can be briefly described as follows:

$$\mathbf{x} = \tilde{\mathbf{h}}_0 \xrightarrow{\mathcal{L}_0} \mathbf{h}_1 \xrightarrow{\sigma} \tilde{\mathbf{h}}_1 \quad \dots \quad \xrightarrow{\mathcal{L}_{L-1}} \mathbf{h}_L \xrightarrow{\sigma} \tilde{\mathbf{h}}_L \xrightarrow{\mathcal{L}_L} \mathbf{h}_{L+1} = \phi(\mathbf{x}), \quad (2.2)$$

where  $N_0 = d \in \mathbb{N}^+$ ,  $N_1, N_2, \dots, N_L \in \mathbb{N}^+$ ,  $N_{L+1} = 1$ ,  $\mathbf{A}_i \in \mathbb{R}^{N_{i+1} \times N_i}$  and  $\mathbf{b}_i \in \mathbb{R}^{N_{i+1}}$  are the weight matrix and the bias vector in the  $i$ -th affine linear transform  $\mathcal{L}_i$ , respectively, i.e.,

$$\mathbf{h}_{i+1} = \mathbf{A}_i \cdot \tilde{\mathbf{h}}_i + \mathbf{b}_i =: \mathcal{L}_i(\tilde{\mathbf{h}}_i) \quad \text{for } i = 0, 1, \dots, L$$

and

$$\tilde{h}_{i,j} = \sigma(h_{i,j}) \quad \text{for } j = 1, 2, \dots, N_i \text{ and } i = 1, 2, \dots, L.$$

Here,  $\tilde{h}_{i,j}$  and  $h_{i,j}$  are the  $j$ -th entry of  $\tilde{\mathbf{h}}_i$  and  $\mathbf{h}_i$ , respectively, for  $j = 1, 2, \dots, N_i$  and  $i = 1, 2, \dots, L$ . If  $\sigma$  is applied to a vector entrywisely, i.e.,

$$\sigma(\mathbf{y}) = (\sigma(y_1), \dots, \sigma(y_d)) \quad \text{for any } \mathbf{y} = (y_1, \dots, y_d) \in \mathbb{R}^d,$$

then  $\phi$  can be represented in a form of function compositions as follows:

$$\phi(\mathbf{x}) = \mathcal{L}_L \circ \sigma \circ \dots \circ \sigma \circ \mathcal{L}_1 \circ \sigma \circ \mathcal{L}_0(\mathbf{x}) \quad \text{for any } \mathbf{x} \in \mathbb{R}^d.$$

See Figure 2 for an example.



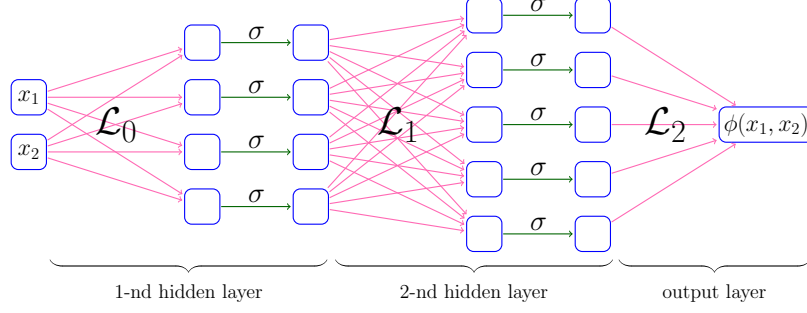


Figure 2: An example of a  $\sigma$ -activated network with width 5 and depth 2.

## 2.2 Proof of Theorem 1.1

To prove Theorems 1.1 and 1.2, we introduce an auxiliary theorem below with a similar result ignoring the approximation inside the trifling region.

**Theorem 2.1.** *Given any  $n \in \mathbb{N}^+$ , there exist  $\phi_1 \in \mathcal{H}_{2^{dn+4}}(d, 1)$  and  $\phi_2 \in \mathcal{H}_{2^{dn+5}n}(n, 1)$  such that: For any continuous function  $f : [0, 1]^d \rightarrow [0, 1]$ , there exists a linear map  $\mathcal{L} : \mathbb{R} \rightarrow \mathbb{R}^n$  satisfying  $\|\phi_2 \circ \mathcal{L} \circ \phi_1\|_{L^\infty(\mathbb{R}^d)} \leq 1$  and*

$$|\phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x}) - f(\mathbf{x})| \leq \omega_f(\sqrt{d}2^{-n}) + 2^{-n} \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta),$$

where  $K = 2^n$ ,  $\delta$  is an arbitrary number in  $(0, \frac{1}{3K}]$ , and  $\mathcal{L}$  is given by  $\mathcal{L}(t) = (a_1 t, a_2 t, \dots, a_n t)$ . Here,  $a_i \in [0, \frac{1}{3})$  is determined by  $f$  and  $n$ , and  $2^m a_i \in \mathbb{N}$  for  $i = 1, 2, \dots, n$ , where  $m = 2^{dn+1}$ .

The proof of Theorem 2.1 can be found later in Section 3. Let us first prove Theorem 1.1 based on Theorem 2.1.

*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  $r > 0$ . Set  $s = 2\omega_f(\sqrt{d}) > 0$  and  $b = f(\mathbf{0}) - \omega_f(\sqrt{d})$ . Then, by defining

$$\tilde{f} := \frac{f - b}{s} = \frac{f - f(\mathbf{0}) + \omega_f(\sqrt{d})}{2\omega_f(\sqrt{d})},$$

we have  $\tilde{f}(\mathbf{x}) \in [0, 1]$  for any  $\mathbf{x} \in [0, 1]^d$ . By applying Theorem 2.1 to  $\tilde{f}$ , there exist two functions,  $\phi_1 : \mathbb{R}^d \rightarrow \mathbb{R}$  and  $\phi_2 : \mathbb{R}^n \rightarrow \mathbb{R}$ , both of which are independent of  $f$  and can be implemented by ReLU networks with  $\leq 2^{dn+4}$  and  $\leq 2^{dn+5}n$  parameters, respectively, and a linear map  $\mathcal{L} : \mathbb{R} \rightarrow \mathbb{R}^n$  satisfying  $\|\phi_2 \circ \mathcal{L} \circ \phi_1\|_{L^\infty(\mathbb{R}^d)} \leq 1$  and

$$|\phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x}) - \tilde{f}(\mathbf{x})| \leq \omega_{\tilde{f}}(\sqrt{d}2^{-n}) + 2^{-n} \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta),$$

where  $K = 2^n$ ,  $\delta$  is an arbitrary number in  $(0, \frac{1}{3K}]$ , and  $\mathcal{L}$  is given by  $\mathcal{L}(t) = (a_1 t, a_2 t, \dots, a_n t)$  with  $a_1, a_2, \dots, a_n \in [0, \frac{1}{3})$  determined by  $\tilde{f}$  and  $n$ . Since  $\tilde{f}$  is derived from  $f$ ,  $a_1, a_2, \dots, a_n$  are essentially determined by  $f$  and  $n$ .

Then choose a small  $\delta$  satisfying

$$dK\delta 2^p \leq 2^{-pn} = (2^{-n})^p.$$

Note that the Lebesgue measure of  $\Omega([0, 1]^d, K, \delta)$  is bounded by  $dK\delta$  and

$$|\phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x}) - \tilde{f}(\mathbf{x})| \leq |\phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x})| + |\tilde{f}(\mathbf{x})| \leq 1 + 1 = 2 \quad \text{for any } \mathbf{x} \in [0, 1]^d.$$

Then,  $\|\phi_2 \circ \mathcal{L} \circ \phi_1 - \tilde{f}\|_{L^p([0, 1]^d)}^p$  is bounded by

$$\begin{aligned} & \int_{[0, 1]^d \setminus \Omega([0, 1]^d, K, \delta)} |\phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x}) - \tilde{f}(\mathbf{x})|^p d\mathbf{x} + \int_{\Omega([0, 1]^d, K, \delta)} |\phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x}) - \tilde{f}(\mathbf{x})|^p d\mathbf{x} \\ & \leq (\omega_{\tilde{f}}(\sqrt{d}2^{-n}) + 2^{-n})^p + dK\delta 2^p \leq (\omega_{\tilde{f}}(\sqrt{d}2^{-n}) + 2^{-n})^p + (2^{-n})^p \leq (\omega_{\tilde{f}}(\sqrt{d}2^{-n}) + 2^{-n+1})^p, \end{aligned}$$

implying  $\|\phi_2 \circ \mathcal{L} \circ \phi_1 - \tilde{f}\|_{L^p([0, 1]^d)} \leq \omega_{\tilde{f}}(\sqrt{d}2^{-n}) + 2^{-n+1}$ . Note that  $\omega_f(r) = s \cdot \omega_{\tilde{f}}(r)$  for any  $r \geq 0$ . Therefore, we have

$$\begin{aligned} & \|s \cdot (\phi_2 \circ \mathcal{L} \circ \phi_1) + b - f\|_{L^p([0, 1]^d)} = \left\| s \cdot (\phi_2 \circ \mathcal{L} \circ \phi_1) + b - (s \cdot \tilde{f} + b) \right\|_{L^p([0, 1]^d)} \\ & = s \|\phi_2 \circ \mathcal{L} \circ \phi_1 - \tilde{f}\|_{L^p([0, 1]^d)} \leq s \cdot \omega_{\tilde{f}}(\sqrt{d}2^{-n}) + 2^{-n+1}s = \omega_f(\sqrt{d}2^{-n}) + 2^{-n+2}\omega_f(\sqrt{d}). \end{aligned}$$

So we finish the proof.  $\square$

### 2.3 Proof of Theorem 1.2

Next, let us prove Theorem 1.2. To this end, we need to introduce the following lemma, which is actually Lemma 3.4 of [25] (or Lemma 3.11 of [48]).

**Lemma 2.2** (Lemma 3.4 of [25]). *Given any  $\varepsilon > 0$ ,  $K \in \mathbb{N}^+$ , and  $\delta \in (0, \frac{1}{3K}]$ , assume  $f \in C([0, 1]^d)$  and  $g: \mathbb{R}^d \rightarrow \mathbb{R}$  is a general function with*

$$|g(\mathbf{x}) - f(\mathbf{x})| \leq \varepsilon \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta).$$

Then

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

where  $\phi := \phi_d$  is defined by induction through

$$\phi_{i+1}(\mathbf{x}) := \text{mid}(\phi_i(\mathbf{x} - \delta \mathbf{e}_{i+1}), \phi_i(\mathbf{x}), \phi_i(\mathbf{x} + \delta \mathbf{e}_{i+1})) \quad \text{for } i = 0, 1, \dots, d-1,$$

where  $\phi_0 = g$  and  $\{\mathbf{e}_i\}_{i=1}^d$  is the standard basis in  $\mathbb{R}^d$ .

With Lemma 2.2 in hand, we are ready to present the proof of Theorem 1.2.

*Proof of Theorem 1.2.* We may assume  $f$  is not a constant function since it is a trivial case. Then  $\omega_f(r) > 0$  for any  $r > 0$ . Set  $s = 2\omega_f(\sqrt{d}) > 0$  and  $b = f(\mathbf{0}) - \omega_f(\sqrt{d})$ . Then, by defining

$$\tilde{f} := \frac{f - b}{s} = \frac{f - f(\mathbf{0}) + \omega_f(\sqrt{d})}{2\omega_f(\sqrt{d})},$$

we have  $\tilde{f}(\mathbf{x}) \in [0, 1]$  for any  $\mathbf{x} \in [0, 1]^d$ . By applying Theorem 2.1 to  $\tilde{f}$ , there exist two functions,  $\psi_{1,0}: \mathbb{R}^d \rightarrow \mathbb{R}$  and  $\psi_{2,0}: \mathbb{R}^n \rightarrow \mathbb{R}$ , both of which are independent of  $\tilde{f}$  (or

f) and can be implemented by ReLU networks with  $\leq 2^{dn+4}$  and  $\leq 2^{dn+5}n$  parameters, respectively, and a linear map  $\mathcal{L}_{\mathbf{a},0} : \mathbb{R} \rightarrow \mathbb{R}^n$  satisfying  $\|\psi_0\|_{L^\infty(\mathbb{R}^d)} \leq 1$  and

$$|\psi_0(\mathbf{x}) - \tilde{f}(\mathbf{x})| \leq \omega_{\tilde{f}}(\sqrt{d}2^{-n}) + 2^{-n} =: \varepsilon \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta),$$

where  $\psi_0 := \psi_{2,0} \circ \mathcal{L}_{\mathbf{a},0} \circ \psi_{1,0}$ ,  $K = 2^n$ ,  $\delta$  is an arbitrary number in  $(0, \frac{1}{3K}]$ , and  $\mathcal{L}_{\mathbf{a},0}$  is given by  $\mathcal{L}_{\mathbf{a},0}(t) = (a_1 t, a_2 t, \dots, a_n t)$  with  $\mathbf{a} = (a_1, a_2, \dots, a_n)$  determined by  $f$  and  $n$ . Moreover,  $a_i \in [0, \frac{1}{3})$  and  $2^m a_i \in \mathbb{N}$  for  $i = 1, 2, \dots, n$ , where  $m = 2^{dn+1}$ .<sup>①</sup>

Choose a small  $\delta$  satisfying  $d \cdot \omega_{\tilde{f}}(\delta) \leq 2^{-n}$ . With  $\psi_0 = \psi_{2,0} \circ \mathcal{L}_{\mathbf{a},0} \circ \psi_{1,0}$  in hand, we can define  $\psi_1, \dots, \psi_d$  by induction via

$$\psi_{i+1}(\mathbf{x}) := \text{mid}\left(\psi_i(\mathbf{x} - \delta \mathbf{e}_{i+1}), \psi_i(\mathbf{x}), \psi_i(\mathbf{x} + \delta \mathbf{e}_{i+1})\right) \quad \text{for } i = 0, 1, \dots, d-1.$$

The detailed iterative equations for  $\psi_i : \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $\boldsymbol{\psi}_{1,i} : \mathbb{R}^d \rightarrow \mathbb{R}^{3^i}$ ,  $\mathcal{L}_{\mathbf{a},i} : \mathbb{R}^{3^i} \rightarrow \mathbb{R}^{3^i d}$ , and  $\psi_{2,i} : \mathbb{R}^{3^i d} \rightarrow \mathbb{R}$ , for  $i = 1, 2, \dots, d$ , are listed as follows.

- $\psi_i = \psi_{2,i} \circ \mathcal{L}_{\mathbf{a},i} \circ \boldsymbol{\psi}_{1,i}$ .
- $\boldsymbol{\psi}_{1,i}(\mathbf{y}) = \left(\boldsymbol{\psi}_{1,i-1}(\mathbf{y} - \delta \mathbf{e}_i), \boldsymbol{\psi}_{1,i-1}(\mathbf{y}), \boldsymbol{\psi}_{1,i-1}(\mathbf{y} + \delta \mathbf{e}_i)\right) \quad \text{for any } \mathbf{y} \in \mathbb{R}^d$ .
- $\mathcal{L}_{\mathbf{a},i}(\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3) = \left(\mathcal{L}_{\mathbf{a},i-1}(\mathbf{y}_1), \mathcal{L}_{\mathbf{a},i-1}(\mathbf{y}_2), \mathcal{L}_{\mathbf{a},i-1}(\mathbf{y}_3)\right) \quad \text{for any } \mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3 \in \mathbb{R}^{3^{i-1}}$ .
- $\psi_{2,i}(\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3) = \text{mid}\left(\psi_{2,i-1}(\mathbf{y}_1), \psi_{2,i-1}(\mathbf{y}_2), \psi_{2,i-1}(\mathbf{y}_3)\right) \quad \text{for any } \mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3 \in \mathbb{R}^{3^{i-1}d}$ .

See the illustrations in Figure 3.

By Lemma 2.2, we have

$$|\phi(\mathbf{x}) - \tilde{f}(\mathbf{x})| \leq \varepsilon + d \cdot \omega_{\tilde{f}}(\delta) \leq \omega_{\tilde{f}}(\sqrt{d}2^{-n}) + 2^{-n+1} \quad \text{for any } \mathbf{x} \in [0, 1]^d,$$

where  $\phi := \psi_d = \psi_{2,d} \circ \mathcal{L}_{\mathbf{a},d} \circ \boldsymbol{\psi}_{1,d}$ . By defining  $\mathcal{L} := \mathcal{L}_{\mathbf{a},d}$ ,  $\boldsymbol{\phi}_1 := \boldsymbol{\psi}_{1,d}$ , and  $\phi_2 := \psi_{2,d}$ , we have

$$|\phi_2 \circ \mathcal{L} \circ \boldsymbol{\phi}_1(\mathbf{x}) - \tilde{f}(\mathbf{x})| \leq \omega_{\tilde{f}}(\sqrt{d}2^{-n}) + 2^{-n+1} \quad \text{for any } \mathbf{x} \in [0, 1]^d,$$

As shown in Figure 3,  $\mathcal{L} := \mathcal{L}_{\mathbf{a},d}$  is a linear map from  $\mathbb{R}^{3^d}$  to  $\mathbb{R}^{3^d n}$  determined by  $\mathbf{a} = (a_1, a_2, \dots, a_n) \in [0, \frac{1}{3})^n$ , which depends on  $f$  and  $n$ . Moreover, as shown in Figure 3,  $\boldsymbol{\phi}_1 := \boldsymbol{\psi}_{1,d}$  and  $\phi_2 := \psi_{2,d}$  are independent of  $f$  and can be implemented by ReLU networks with

$$\leq 3^d(2^{dn+4}) + 3d(d+1)(3^{d-1} + 3^{d-2} + \dots + 3^0) \leq 3^d 2^{dn+5}$$

and

$$\leq 3^d(2^{dn+5}n) + 280(3^{d-1} + 3^{d-2} + \dots + 3^0) \leq 3^d 2^{dn+8}n \quad \text{②}$$

parameters, respectively.

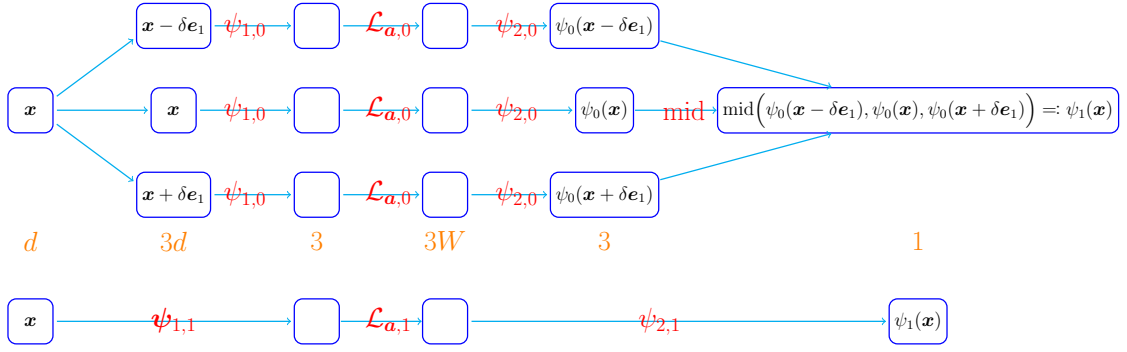
Note that  $\omega_f(r) = s \cdot \omega_{\tilde{f}}(r)$  for any  $r \geq 0$ . Therefore, we have

$$\begin{aligned} & \|s \cdot (\phi_2 \circ \mathcal{L} \circ \boldsymbol{\phi}_1) + b - f\|_{L^\infty([0,1]^d)} = \left\| s \cdot (\phi_2 \circ \mathcal{L} \circ \boldsymbol{\phi}_1) + b - (s \cdot \tilde{f} + b) \right\|_{L^\infty([0,1]^d)} \\ & = s \|\phi_2 \circ \mathcal{L} \circ \boldsymbol{\phi}_1 - \tilde{f}\|_{L^\infty([0,1]^d)} \leq s \cdot \omega_{\tilde{f}}(\sqrt{d}2^{-n}) + 2^{-n+1}s = \omega_f(\sqrt{d}2^{-n}) + 2^{-n+2}\omega_f(\sqrt{d}). \end{aligned}$$

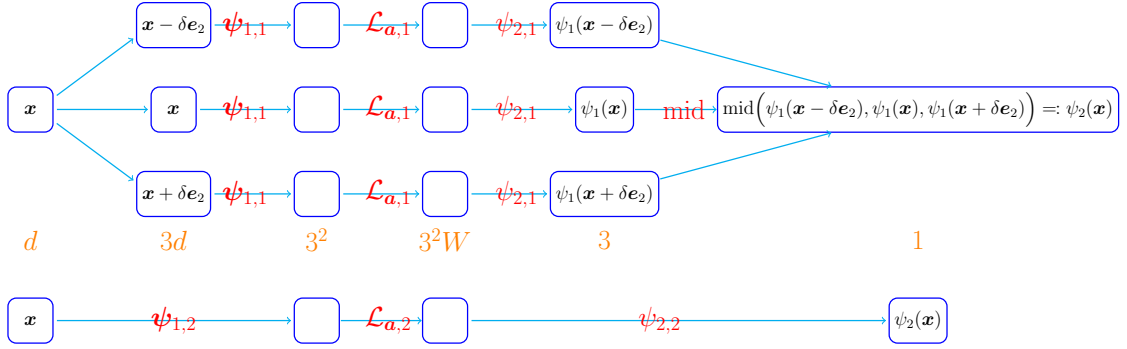
So we finish the proof.  $\square$

<sup>①</sup>This property will be used in the proof of Theorem 1.3.

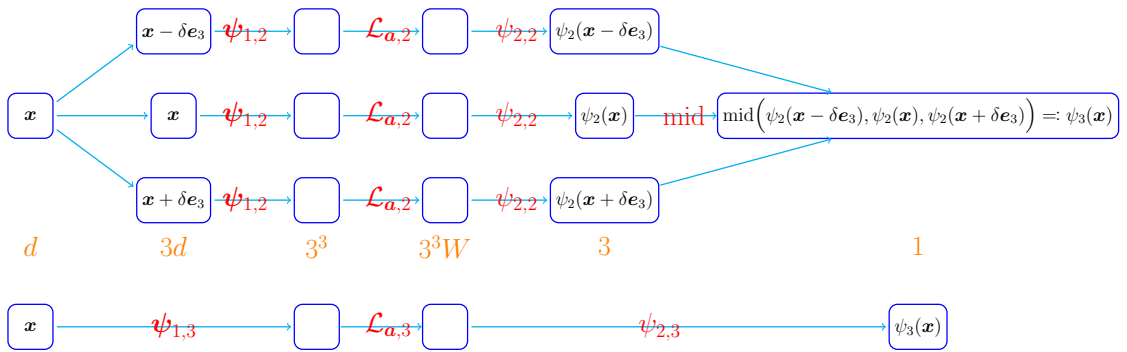
<sup>②</sup>As shown Lemma 3.1 of [25], “mid( $\cdot, \cdot, \cdot$ )” can be implemented by a ReLU network with width 14 and depth 2, which has  $\leq (3+1) \times 14 + (14+1) \times 14 + (14+1) = 280$  parameters.



(a) An illustration of the network architecture implementing  $\psi_1 = \psi_{2,1} \circ \mathcal{L}_{a,1} \circ \psi_{1,1}$  based on  $\psi_0 = \psi_{2,0} \circ \mathcal{L}_{a,0} \circ \psi_{1,0}$ . The top architecture is in detail, while the bottom one is just a sketch of the top one. The orange numbers indicate the number of neurons in each layer.



(b) An illustration of the network architecture implementing  $\psi_2 = \psi_{2,2} \circ \mathcal{L}_{a,2} \circ \psi_{1,2}$  based on  $\psi_1 = \psi_{2,1} \circ \mathcal{L}_{a,1} \circ \psi_{1,1}$ . The top architecture is in detail, while the bottom one is just a sketch of the top one. The orange numbers indicate the number of neurons in each layer.



(c) An illustration of the network architecture implementing  $\psi_3 = \psi_{2,3} \circ \mathcal{L}_{a,3} \circ \psi_{1,3}$  based on  $\psi_2 = \psi_{2,2} \circ \mathcal{L}_{a,2} \circ \psi_{1,2}$ . The top architecture is in detail, while the bottom one is just a sketch of the top one. The orange numbers indicate the number of neurons in each layer.

Figure 3: Illustrations of the implementations of  $\phi_1$ ,  $\phi_2$ , and  $\phi_3$ . The inductive implementations of  $\phi_4, \dots, \phi_d$  are similar.

## 2.4 Proof of Theorem 1.3

To simplify the proof of Theorem 1.3, we introduce two lemmas below. First, we need to establish a lemma showing how to store many parameters in one intrinsic parameter via a fixed network.

**Lemma 2.3.** *Given any  $m, n \in \mathbb{N}$ , there exists a vector-valued function  $\phi : \mathbb{R} \rightarrow \mathbb{R}^n$  realized by a ReLU network such that: For any  $a_i \in [0, 1)$  with  $2^m a_i \in \mathbb{N}$  for  $i = 0, 1, \dots, n$ , there exists a real number  $v \in [0, 1)$  such that*

$$\phi(v) = (a_1, a_2, \dots, a_n).$$

Next, we establish another lemma using a ReLU network to uniformly approximate multiplication operation  $\psi(x, y) = xy$  well.

**Lemma 2.4.** *For any  $M > 0$  and  $\eta > 0$ , there exists a function  $\psi_\eta : \mathbb{R}^2 \rightarrow \mathbb{R}$  realized by a ReLU network such that*

$$\psi_\eta(x, y) \rightrightarrows \psi(x, y) = xy \quad \text{on } [-M, M]^2 \quad \text{as } \eta \rightarrow 0^+,$$

where  $\rightrightarrows$  denotes the uniform convergence.

The proof of Lemma 2.3 is placed later in this section. Lemma 2.4 is just a direct result of Lemma 4.2 of [25]. With Lemmas 2.3 and 2.4 in hand, we are ready to prove Theorem 1.3.

*Proof of Theorem 1.3.* For any  $\varepsilon > 0$ , choose a large  $n = n(\varepsilon, \alpha, \lambda) \in \mathbb{N}^+$  such that

$$5\lambda d^{\alpha/2} 2^{-\alpha n} \leq \varepsilon/2.$$

Since  $f$  is a Hölder continuous function on  $[0, 1]^d$  of order  $\alpha \in (0, 1]$  with a Hölder constant  $\lambda > 0$ , we have  $\omega_f(r) \leq \lambda r^\alpha$  for any  $r \geq 0$ . By Theorem 1.2, there exist two functions  $\phi_1 : \mathbb{R}^d \rightarrow \mathbb{R}^{3^d}$  and  $\phi_2 : \mathbb{R}^{3^d n} \rightarrow \mathbb{R}$ , implemented by  $f$ -independent ReLU networks, such that

$$\|s(\phi_2 \circ \mathcal{L} \circ \phi_1) + b - f\|_{L^\infty([0, 1]^d)} \leq \omega_f(\sqrt{d} 2^{-n}) + 2^{-n+2} \omega_f(\sqrt{d}) \leq 5\lambda d^{\alpha/2} 2^{-\alpha n} \leq \varepsilon/2, \quad (2.3)$$

where  $s = 2\omega_f(\sqrt{d}) \leq 2\lambda d^{\alpha/2}$ ,  $b = f(\mathbf{0}) - \omega_f(\sqrt{d})$ , and  $\mathcal{L} : \mathbb{R}^{3^d} \rightarrow \mathbb{R}^{3^d n}$  is a linear map given by

$$\mathcal{L}(y_1, \dots, y_{3^d}) = \left( \mathcal{L}_0(y_1), \dots, \mathcal{L}_0(y_{3^d}) \right) \quad \text{for any } \mathbf{y} = (y_1, \dots, y_{3^d}) \in \mathbb{R}^{3^d},$$

where  $\mathcal{L}_0 : \mathbb{R} \rightarrow \mathbb{R}^n$  is a linear map given by  $\mathcal{L}_0(t) = (a_1 t, a_2 t, \dots, a_n t)$  with  $a_1, a_2, \dots, a_n \in [0, \frac{1}{3})$  determined by  $f$  and  $n$ . Since  $a_1, a_2, \dots, a_n$  are repeated  $3^d$  times in the definition of  $\mathcal{L}$ , there are  $3^d n$  parameters in total. We will show how to store these  $3^d n$  parameters in one intrinsic parameter  $v$  via an  $f$ -independent ReLU network as shown in the following two steps.

- Regard  $a_1, a_2, \dots, a_n$  as inputs, but not parameters. See the difference in Figure 4. Then, we only need to store  $a_1, a_2, \dots, a_n$  one time by copying them  $3^d$  times.

- As stated in the proof of Theorem 1.2,  $a_1, a_2, \dots, a_n$  have finite binary representations. Then, we can store them in a key parameter  $v$  and use an  $f$ -independent ReLU network to extract them from  $v$ .

The details of these two steps can be found below.

**Step 1:** Regard  $a_1, a_2, \dots, a_n$  as inputs and copy them  $3^d$  times.

Since  $a_1, a_2, \dots, a_n$  are regarded as inputs, the implementation of  $\mathcal{L}_0(t) = (a_1 t, a_2 t, \dots, a_n t)$  requires multiplication operations. This means that we need to approximate  $\psi(x, y) = xy$  well via an  $f$ -independent ReLU network. See Figure 4 for illustrations.

Denote  $\mathbf{a} = (a_1, a_2, \dots, a_n)$  and

$$\phi_1(\mathbf{x}) = \left( \phi_{1,1}(\mathbf{x}), \phi_{1,2}(\mathbf{x}), \dots, \phi_{1,3^d}(\mathbf{x}) \right) \quad \text{for any } \mathbf{x} \in [0, 1]^d.$$

Then define

$$M := 1 + \sup \left\{ |\phi_{1,j}(\mathbf{x})| : \mathbf{x} \in [0, 1]^d, j = 1, 2, \dots, 3^d \right\}.$$

By Lemma 2.4, there exists a function  $\psi_\eta : \mathbb{R}^2 \rightarrow \mathbb{R}$  realized by a ReLU network such that

$$\psi_\eta(x, y) \rightrightarrows \psi(x, y) = xy \quad \text{on } [-M, M]^2 \quad \text{as } \eta \rightarrow 0^+.$$

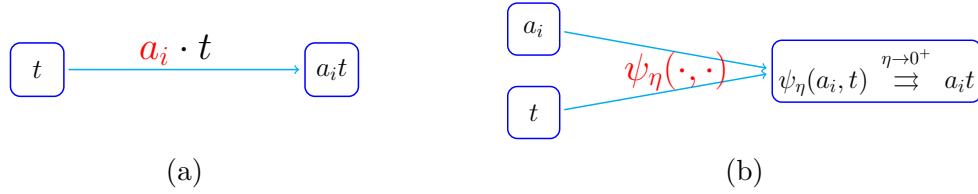


Figure 4: Illustrations of two methods getting/approximating  $a_i t$  for  $i = 1, 2, \dots, n$  and  $t \in \{\phi_{1,j}(\mathbf{x}) : \mathbf{x} \in [0, 1]^d, j = 1, 2, \dots, 3^d\}$ . (a) By regarding  $a_i$  as a parameter, one can easily get the product of an input  $t$  and an parameter  $a_i$ . (b) By regarding  $a_i$  as an input, one needs to use a ReLU network to approximate the multiplication operation for approximating  $a_i t$  well.

Note that  $\mathcal{L}$  and  $\mathcal{L}_0$  depend on  $\mathbf{a} = (a_1, a_2, \dots, a_n)$ . For clarity, we denote  $\mathcal{L}_{\mathbf{a}} = \mathcal{L}$  and  $\mathcal{L}_{\mathbf{a},0} = \mathcal{L}_0$ . For any  $a_1, a_2, \dots, a_n \in [0, \frac{1}{3}] \subseteq [-M, M]$  and

$$t \in \left\{ \phi_{1,j}(\mathbf{x}) : \mathbf{x} \in [0, 1]^d, j = 1, 2, \dots, 3^d \right\} \subseteq [-M, M],$$

we can use

$$\mathcal{L}_{\mathbf{a},0,\eta}(t) := \left( \psi_\eta(a_1, t), \psi_\eta(a_2, t), \dots, \psi_\eta(a_n, t) \right)$$

to approximate

$$\mathcal{L}_{\mathbf{a},0}(t) = (a_1 t, a_2 t, \dots, a_n t) = \left( \psi(a_1, t), \psi(a_2, t), \dots, \psi(a_n, t) \right).$$

Then

$$\mathcal{L}_{\mathbf{a},\eta}(y_1, \dots, y_{3^d}) := \left( \mathcal{L}_{\mathbf{a},0,\eta}(y_1), \dots, \mathcal{L}_{\mathbf{a},0,\eta}(y_{3^d}) \right)$$

can also approximate

$$\mathcal{L}_{\mathbf{a}}(y_1, \dots, y_{3^d}) = \left( \mathcal{L}_{\mathbf{a},0}(y_1), \dots, \mathcal{L}_{\mathbf{a},0}(y_{3^d}) \right)$$

well. Define  $\tilde{\phi}_\eta : \mathbb{R}^{n+3^d} \rightarrow \mathbb{R}^{3^d n}$  and  $\tilde{\phi} : \mathbb{R}^{n+3^d} \rightarrow \mathbb{R}^{3^d n}$  via

$$\tilde{\phi}_\eta(\mathbf{y}, \mathbf{a}) := \mathcal{L}_{\mathbf{a},\eta}(y_1, \dots, y_{3^d}) \quad \text{for any } \mathbf{a} \in [0, \frac{1}{3}]^n \text{ and } \mathbf{y} = (y_1, \dots, y_{3^d}) \in [-M, M]^{3^d}$$

and

$$\tilde{\phi}(\mathbf{y}, \mathbf{a}) := \mathcal{L}_{\mathbf{a}}(y_1, \dots, y_{3^d}) \quad \text{for any } \mathbf{a} \in [0, \frac{1}{3}]^n \text{ and } \mathbf{y} = (y_1, \dots, y_{3^d}) \in [-M, M]^{3^d}.$$

Since

$$\psi_\eta(x, y) \rightrightarrows \psi(x, y) = xy \quad \text{on } [-M, M]^2 \quad \text{as } \eta \rightarrow 0^+,$$

it is easy to verify that

$$\tilde{\phi}_\eta(\mathbf{y}, \mathbf{a}) \rightrightarrows \tilde{\phi}(\mathbf{y}, \mathbf{a}) \quad \text{for } \mathbf{y} \in [-M, M]^{3^d} \text{ and } \mathbf{a} \in [0, \frac{1}{3}]^n \quad \text{as } \eta \rightarrow 0^+.$$

Note that  $\phi_1(\mathbf{x}) \in [-M, M]^{3^d}$  for any  $\mathbf{x} \in [0, 1]^d$ . Then,

$$\phi_2 \circ \tilde{\phi}_\eta(\phi_1(\mathbf{x}), \mathbf{a}) \rightrightarrows \phi_2 \circ \tilde{\phi}(\phi_1(\mathbf{x}), \mathbf{a}) \quad \text{for } \mathbf{x} \in [0, 1]^d \text{ and } \mathbf{a} \in [0, \frac{1}{3}]^n \quad \text{as } \eta \rightarrow 0^+.$$

The fact  $\tilde{\phi}(\mathbf{y}, \mathbf{a}) = \mathcal{L}_{\mathbf{a}}(\mathbf{y}) = \mathcal{L}(\mathbf{y})$  implies  $\phi_2 \circ \tilde{\phi}(\phi_1(\mathbf{x}), \mathbf{a}) = \phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x})$ . Therefore,

$$\phi_2 \circ \tilde{\phi}_\eta(\phi_1(\mathbf{x}), \mathbf{a}) \rightrightarrows \phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x}) \quad \text{for } \mathbf{x} \in [0, 1]^d \text{ and } \mathbf{a} \in [0, \frac{1}{3}]^n \quad \text{as } \eta \rightarrow 0^+.$$

Choose a small  $\eta = \eta(n) > 0$  such that

$$\left| \phi_2 \circ \tilde{\phi}_\eta(\phi_1(\mathbf{x}), \mathbf{a}) - \phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x}) \right| \leq 2^{-n} \quad \text{for any } \mathbf{x} \in [0, 1]^d \text{ and } \mathbf{a} \in [0, \frac{1}{3}]^n. \quad (2.4)$$

Recall that  $\psi_\eta$  can be realized by an  $f$ -independent ReLU network. It is easy to verify that  $\tilde{\phi}_\eta$  can also be realized by an  $f$ -independent ReLU network.

**Step 2:** Store  $a_1, a_2, \dots, a_n$  in a key parameter  $v$ .

As we can see from the proof of Theorem 1.2,  $a_i \in [0, \frac{1}{3}]$  with  $2^m a_i \in \mathbb{N}$  for  $i = 1, 2, \dots, n$ , where  $m = 2^{dn+1}$ . That is,

$$a_i \in \left\{ \text{bin}0.\theta_1 \dots \theta_m : \theta_\ell \in \{0, 1\}, \ell = 1, 2, \dots, m \right\}.$$

Then, by Lemma 2.3, there exists a real number  $v \in [0, 1)$  and a vector function  $\phi_0 : \mathbb{R} \rightarrow \mathbb{R}^n$  implemented by a ReLU network independent of  $a_1, a_2, \dots, a_n$  such that

$$\phi_0(v) = (a_1, a_2, \dots, a_n) = \mathbf{a}.$$

Next, we can define

$$\phi(\mathbf{x}, v) := \phi_2 \circ \tilde{\phi}_\eta(\phi_1(\mathbf{x}), \phi_0(v)) = \phi_2 \circ \tilde{\phi}_\eta(\phi_1(\mathbf{x}), \mathbf{a})$$



for any  $\mathbf{x} \in [0, 1]^d$  and  $v \in [0, 1]$ .

Since the ReLU network realizing  $\phi_0$  is independent of  $a_1, a_2, \dots, a_n$ , and hence independent of  $f$ . Recall that  $\phi_1, \phi_2$ , and  $\tilde{\phi}_\eta$  can be implemented by  $f$ -independent ReLU networks. Hence,

$$\phi(\mathbf{x}, v) = \phi_2 \circ \tilde{\phi}_\eta(\phi_1(\mathbf{x}), \phi_0(v))$$

can also be implemented by an  $f$ -independent ReLU network. It remains to estimate the error. By Equations (2.3) and (2.4), we have

$$\begin{aligned} |s\phi(\mathbf{x}, v) + b - f(\mathbf{x})| &\leq \left| s\phi(\mathbf{x}, v) + b - (s\phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x}) + b) \right| + \left| s\phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x}) + b - f(\mathbf{x}) \right| \\ &\leq s|\phi(\mathbf{x}, v) - \phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x})| + \varepsilon/2 \\ &\leq 2\lambda d^{\alpha/2} \left| \phi_2 \circ \tilde{\phi}_\eta(\phi_1(\mathbf{x}), \mathbf{a}) - \phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x}) \right| + \varepsilon/2 \\ &\leq 2\lambda d^{\alpha/2} 2^{-n} + \varepsilon/2 \leq 5\lambda d^{\alpha/2} 2^{-\alpha n} + \varepsilon/2 \leq \varepsilon/2 + \varepsilon/2 = \varepsilon. \end{aligned}$$

So we finish the proof.  $\square$

Finally, let us prove Lemma 2.3 to end this section.

*Proof of Lemma 2.3.* Since  $a_i \in [0, 1]$  with  $2^m a_i \in \mathbb{N}$  for  $i = 1, 2, \dots, n$ ,  $a_i$  can be represented as a binary form

$$a_i = \text{bin}0.a_{i,1}a_{i,2}\dots a_{i,m}.$$

Denote

$$v = \sum_{i=1}^n 2^{-m(i-1)} a_i = \text{bin}0.\underbrace{a_{1,1}\dots a_{1,m}}_{\text{store } a_1} \underbrace{a_{2,1}\dots a_{2,m}}_{\text{store } a_2} \dots \underbrace{a_{n,1}\dots a_{n,m}}_{\text{store } a_n},$$

which requires pretty high precision. It is easy to extract  $a_i$  from  $v$  via the floor function ( $\lfloor \cdot \rfloor$ ), i.e.,

$$a_i = \lfloor 2^{mi}v \rfloor / 2^m - \lfloor 2^{m(i-1)}v \rfloor \quad \text{for } i = 1, 2, \dots, n.$$

Next, we need to use a ReLU network to replace the floor function. Let  $g : \mathbb{R} \rightarrow \mathbb{R}$  be the continuous piecewise linear function with the following breakpoints:

$$(\ell, \ell) \quad \text{and} \quad (\ell + 1 - \delta, \ell) \quad \text{for } \ell = 0, 1, \dots, 2^{mn} - 1, \quad \text{where } \delta = 2^{-mn}.$$

Clearly,  $g$  can be realized by a ReLU network independent of  $a_1, a_2, \dots, a_n$ . By defining

$$g_i(t) := g(2^{mi}t)/2^m - g(2^{m(i-1)}t) \quad \text{for } i = 1, 2, \dots, n \text{ and any } t \in \mathbb{R},$$

we have

$$a_i = g(2^{mi}v)/2^m - g(2^{m(i-1)}v) = g_i(v) \quad \text{for } i = 1, 2, \dots, n.$$

Next, The target function  $\phi$  can be defined via

$$\phi(t) = \left( g_1(t), g_2(t), \dots, g_n(t) \right) \quad \text{for any } t \in \mathbb{R}.$$

Thus, we have

$$\phi(v) = \left( g_1(v), g_2(v), \dots, g_n(v) \right) = (a_1, a_2, \dots, a_n).$$

and  $\phi$  can be realized by a ReLU network independent of  $a_1, a_2, \dots, a_n$ . So we finish the proof.  $\square$

### 3 Proof of Theorem 2.1

In this section, we first present the proof sketch of Theorem 2.1 in Section 3.1, and then give the detailed proof in Section 3.2 based on Proposition 3.1, which will be proved later in Section 3.3.

#### 3.1 Sketch of proof

Before proving Theorem 2.1, let us present the key steps as follows.

1. Set  $K = 2^n$ , divide  $[0, 1]^d$  into  $K^d$  cubes  $Q_\beta$  for  $\beta \in \{0, 1, \dots, K-1\}^d$  and the trifling region  $\Omega([0, 1]^d, K, \delta)$ , and denote  $\mathbf{x}_\beta$  as the vertex of  $Q_\beta$  with minimum  $\|\cdot\|_1$  norm for each  $\beta$ . See Figure 5 for illustrations.
2. Design a ReLU sub-network to implement a function  $\phi_1 : \mathbb{R}^d \rightarrow \mathbb{R}$ , independent of  $f$ , projecting the whole  $Q_\beta$  to a number determined by  $\beta$  in  $\{4^j : j = 1, 2, \dots, K^d\}$  for each  $\beta$ .
3. Design a linear map  $\mathcal{L} : \mathbb{R} \rightarrow \mathbb{R}^n$ , given by  $\mathcal{L}(t) = (a_1 t, a_2 t, \dots, a_n t)$ , for later use, where  $a_1, a_2, \dots, a_n$  are determined by  $f$  and  $n$ .
4. Design a ReLU sub-network to implement a function  $\phi_2 : \mathbb{R}^n \rightarrow \mathbb{R}$ , independent of  $f$ , such that  $\phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x}) \approx f(\mathbf{x}_\beta)$  for any  $\mathbf{x} \in Q_\beta$  and each  $\beta \in \{0, 1, \dots, K-1\}^d$ . Then  $\phi_2 \circ \mathcal{L} \circ \phi_1$  approximates  $f$  well outside of  $\Omega([0, 1]^d, K, \delta)$ .
5. Estimate the approximation error of  $\phi_2 \circ \mathcal{L} \circ \phi_1 \approx f$ .

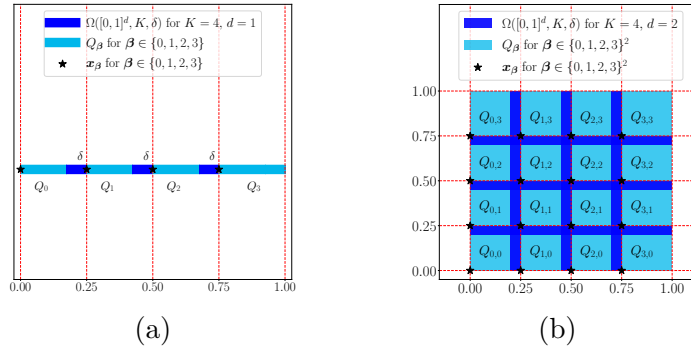


Figure 5: Illustrations of  $\Omega([0, 1]^d, K, \delta)$ ,  $Q_\beta$  and  $\mathbf{x}_\beta$  for  $\beta \in \{0, 1, \dots, K-1\}^d$ . (a)  $K=4, d=1$ . (b)  $K=4, d=2$ .

As we shall see later, the constructions of  $\phi_1$  and  $\mathcal{L}$  are not difficult. The most technical part is to design  $\phi_2$  implemented by a ReLU network, which relies on the following proposition.

**Proposition 3.1.** *Given any  $J \in \mathbb{N}^+$ , there exists a function  $\phi$  implemented by a ReLU network with width 2 and depth  $2J+2$  such that: For any  $\theta_1, \theta_2, \dots, \theta_J \in \{0, 1\}$ , we have  $\phi(x) \in [0, 1]$  for any  $x \in \mathbb{R}$  and*

$$\phi(4^j a) = \theta_j \quad \text{for } j = 1, 2, \dots, J, \quad \text{where } a = \sum_{j=1}^J \theta_j 4^{-j}. \quad (3.1)$$

The proof of this proposition can be found in Section 3.3. We shall point out that the function  $\phi$  in this proposition is independent of  $\theta_1, \theta_2, \dots, \theta_J \in \{0, 1\}$ .

### 3.2 Constructive proof

Now we are ready to give the detailed proof of Theorem 2.1.

*Proof of Theorem 2.1.* The proof consists of five steps.

**Step 1:** Set up.

Set  $K = 2^n$  and let  $\delta > 0$  be a small number determined later. Then define  $\mathbf{x}_\beta := \beta/K$  and divide  $[0, 1]^d$  into  $K^d$  cubes  $Q_\beta$  for  $\beta \in \{0, 1, \dots, K-1\}^d$  and a small region  $\Omega([0, 1]^d, K, \delta)$ . Namely,

$$Q_\beta := \left\{ \mathbf{x} = (x_1, x_2, \dots, x_d) : x_i \in \left[ \frac{\beta_i}{K}, \frac{\beta_i+1}{K} - \delta \right] \text{ for } i = 1, \dots, d \right\},$$

for  $\beta = (\beta_1, \beta_2, \dots, \beta_d) \in \{0, 1, \dots, K-1\}^d$ . Clearly,  $\Omega([0, 1]^d, K, \delta) = [0, 1]^d \setminus (\cup_{\beta \in \{0, 1, \dots, K-1\}^d} Q_\beta)$ . See Figure 5 for illustrations.

**Step 2:** Construct  $\phi_1$ .

Let  $g$  be a “step function” such that

- $g(\frac{k}{K}) = g(\frac{k+1}{K} - \delta) = k$  for  $k = 0, 1, \dots, K-1$  and  $g(1) = K-1$ .
- $g$  is linear between any two adjacent points of

$$\left\{ \frac{k}{K} : k = 0, 1, \dots, K \right\} \cup \left\{ \frac{k+1}{K} - \delta : k = 0, 1, \dots, K-1 \right\}.$$

Then, for any  $\mathbf{x} = (x_1, \dots, x_d) \in Q_\beta$  and  $\beta = (\beta_1, \dots, \beta_d) \in \{0, 1, \dots, K-1\}^d$ , we have

$$g(x_i) = \beta_i \quad \text{for } i = 1, 2, \dots, d.$$

Also, such a function  $g$  can be easily realized by a one-hidden-layer ReLU network with width  $2K$ .

Let  $h$  be a function satisfying  $h(j) = 4^j$  for  $j = 1, 2, \dots, K^d$ . Such a function  $h$  can be easily realized by a one-hidden-layer ReLU network with width  $K^d$ . Then the desired function  $\phi_1 : \mathbb{R}^d \rightarrow \mathbb{R}$  can be defined via

$$\phi_1(\mathbf{x}) := h\left(1 + \sum_{i=1}^d g(x_i)K^{i-1}\right) = h \circ \varrho\left(g(x_1), \dots, g(x_d)\right) \quad \text{for any } \mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d,$$

where  $\varrho : \mathbb{R}^d \rightarrow \mathbb{R}$  is a linear function defined by

$$\varrho(\mathbf{y}) = 1 + \sum_{i=1}^d y_i K^{i-1} \quad \text{for any } \mathbf{y} = (y_1, \dots, y_d) \in \mathbb{R}^d.$$

Clearly,  $\varrho$  is a bijection (one-to-one map) from  $\beta \in \{0, 1, \dots, K-1\}^d$  to  $\varrho(\beta) = 1 + \sum_{i=1}^d \beta_i K^{i-1} \in \{1, 2, \dots, K^d\}$ .

Then, for any  $\mathbf{x} \in Q_\beta$  and  $\beta \in \{0, 1, \dots, K-1\}^d$ , we have

$$\phi_1(\mathbf{x}) = h \circ \varrho(g(x_1), \dots, g(x_d)) = h \circ \varrho(\beta_1, \dots, \beta_d) = 4^{\varrho(\beta)} = 4^{1 + \sum_{i=1}^d \beta_i K^{i-1}}. \quad (3.2)$$

Apparently,  $\phi_1$  is independent of  $f$  and it can be realized by a ReLU network with

$$\leq d(2K \times 3 + 1) + (K + 1) + 3K^d + 1 \leq 11K^d + 2 = 11 \times 2^{dn} + 2 \leq 2^{dn+4}$$

parameters.

**Step 3:** Construct  $\mathcal{L}$ .

For each  $\beta \in \{0, 1, \dots, K-1\}^d$ , it follows from  $f(\mathbf{x}_\beta) \in [0, 1]$  that there exist  $\xi_{\beta,1}, \dots, \xi_{\beta,n}$  such that

$$|f(\mathbf{x}_\beta) - \text{bin}_{0,\xi_{\beta,1}\cdots\xi_{\beta,n}}| \leq 2^{-n}. \quad (3.3)$$

Given any  $j \in \{1, 2, \dots, K^d\}$ , there exists a unique  $\beta \in \{0, 1, \dots, K-1\}^d$  such that  $j = 1 + \sum_{i=1}^d \beta_i K^{i-1} = \varrho(\beta)$ . Thus, for any  $\ell \in \{1, 2, \dots, n\}$ , we can define

$$\theta_{j,\ell} := \xi_{\beta,\ell}, \quad \text{for } j = \varrho(\beta) \text{ and } \beta \in \{0, 1, \dots, K-1\}^d. \quad (3.4)$$

Then the desired linear map  $\mathcal{L}$  can be defined via

$$\mathcal{L}(t) := (a_1 t, a_2 t, \dots, a_n t) \quad \text{for any } t \in \mathbb{R},$$

where  $a_\ell = \sum_{j=1}^{K^d} \theta_{j,\ell} 4^{-j}$  for  $\ell = 1, 2, \dots, n$ . Clearly, for  $\ell = 1, 2, \dots, n$ , we have

$$a_\ell = \sum_{j=1}^{K^d} \theta_{j,\ell} 4^{-j} \in [0, \frac{1}{3}]$$

and

$$2^m a_\ell = 2^{2^{dn+1}} a_\ell = 4^{2^{dn}} a_\ell = 4^{K^d} a_\ell \in \mathbb{N}, \quad \text{where } m = 2^{dn+1}.$$

**Step 4:** Construct  $\phi_2$ .

Fix  $\ell \in \{1, 2, \dots, n\}$ , by Proposition 3.1 (set  $J = K^d$  and  $\theta_j = \theta_{j,\ell}$  therein), there exists a function  $\phi_{2,\ell}$  implemented by a ReLU network with width 2 and depth  $2K^d + 2$  such that  $\phi_{2,\ell}(t) \in [0, 1]$  for any  $t \in \mathbb{R}$  and

$$\phi_{2,\ell}(4^j a_\ell) = \theta_{j,\ell} \quad \text{for } j = 1, 2, \dots, K^d, \quad \text{where } a_\ell = \sum_{j=1}^{K^d} \theta_{j,\ell} 4^{-j}. \quad (3.5)$$

Note that  $\phi_{2,\ell}$  is independent of  $\theta_{j,\ell}$  for  $j = 1, 2, \dots, K^d$ , so it is also independent of  $f$ . Then the desired function  $\phi_2: \mathbb{R}^n \rightarrow \mathbb{R}$  can be defined via

$$\phi_2(\mathbf{y}) := \sum_{\ell=1}^n 2^{-\ell} \phi_{2,\ell}(y_\ell) \quad \text{for any } \mathbf{y} = (y_1, \dots, y_n) \in \mathbb{R}^n.$$

Then  $\phi_2(\mathbf{y}) \in [0, 1]$  for any  $\mathbf{y} \in \mathbb{R}^n$  and  $\phi_2$  can be implemented by a ReLU network, independent of  $f$ , with

$$\leq \left(6(2K^n + 2) + (2 + 1)\right)n + n + 1 = 6n(2^{nd+1} + 2) + 4n + 1 \leq 2^{nd+5}n$$

parameters.

**Step 5:** Estimate the approximation error.

It remains to estimate the approximation error. By Equations (3.2), (3.4), and (3.5), for  $\ell \in \{1, 2, \dots, n\}$ ,  $\mathbf{x} \in Q_\beta$ ,  $j = \varrho(\beta) = 1 + \sum_{i=1}^d \beta_i K^{i-1} \in \{1, 2, \dots, K^d\}$ , and  $\beta \in \{0, 1, \dots, K-1\}^d$ , we have

$$\begin{aligned} \phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x}) &= \phi_2 \circ \mathcal{L} \circ h \circ \varrho(g(x_1), \dots, g(x_d)) = \phi_2 \circ \mathcal{L} \circ h \circ \varrho(\beta) \\ &= \phi_2 \circ \mathcal{L}(4^{\varrho(\beta)}) = \phi_2 \circ \mathcal{L}(4^j) = \phi_2(4^j a_1, 4^j a_2, \dots, 4^j a_n) \\ &= \sum_{\ell=1}^n 2^{-\ell} \phi_{2,\ell}(4^j a_\ell) = \sum_{\ell=1}^n 2^{-\ell} \theta_{j,\ell} = \sum_{\ell=1}^n 2^{-\ell} \xi_{\beta,\ell}. \end{aligned}$$

Then by Equation (3.3), for any  $\mathbf{x} \in Q_\beta$  and  $\beta \in \{0, 1, \dots, K-1\}^d$ , we get

$$\begin{aligned} |f(\mathbf{x}) - \phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x})| &= |f(\mathbf{x}) - f(\mathbf{x}_\beta)| + |f(\mathbf{x}_\beta) - \phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x})| \\ &\leq \omega_f\left(\frac{\sqrt{d}}{K}\right) + |f(\mathbf{x}_\beta) - \sum_{\ell=1}^n 2^{-\ell} \xi_{\beta,\ell}| \\ &\leq \omega_f\left(\frac{\sqrt{d}}{K}\right) + |f(\mathbf{x}_\beta) - \text{bin}0.\xi_{\beta,1}\dots\xi_{\beta,n}| \leq \omega_f(\sqrt{d}2^{-n}) + 2^{-n}. \end{aligned}$$

That is,

$$|f(\mathbf{x}) - \phi_2 \circ \mathcal{L} \circ \phi_1(\mathbf{x})| \leq \omega_f(\sqrt{d}2^{-n}) + 2^{-n} \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta).$$

Moreover, the fact  $\phi_2(\mathbf{y}) \in [0, 1]$  for any  $\mathbf{y} \in \mathbb{R}^n$  implies  $\|\phi_2 \circ \mathcal{L} \circ \phi_1\|_{L^\infty(\mathbb{R}^d)} \leq 1$ . So we finish the proof.  $\square$

### 3.3 Proof of Proposition 3.1

Before proving Proposition 3.1, let us introduce a notation to simplify the proof. We use  $\mathcal{T}_m$  for  $m \in \mathbb{N}^+$  to denote a “sawtooth” function satisfying the following conditions.

- $\mathcal{T}_m : [0, 2m] \rightarrow [0, 1]$  is linear between any two adjacent integers of  $\{0, 1, \dots, 2m\}$ .
- $\mathcal{T}_m(2j) = 0$  for  $j = 0, 1, \dots, m$  and  $\mathcal{T}_m(2j+1) = 1$  for  $j = 0, 1, \dots, m-1$ .

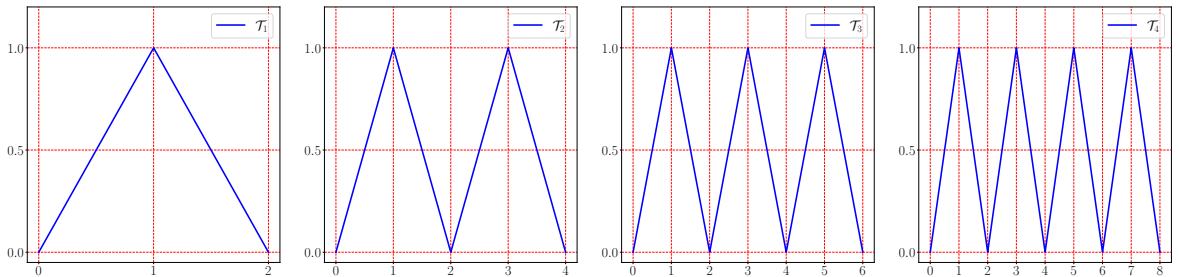


Figure 6: Illustrations of teeth functions  $\mathcal{T}_1$ ,  $\mathcal{T}_2$ ,  $\mathcal{T}_3$  and  $\mathcal{T}_4$ .

To simplify the proof of Proposition 3.1, we first introduce a lemma based on the “sawtooth” function.

**Lemma 3.2.** *Given any  $J \in \mathbb{N}^+$  and  $\theta_j \in \{0, 1\}$  for  $j = 1, 2, \dots, J$ , set  $a = \sum_{j=1}^J \theta_j 4^{-j}$ . Then*

$$\mathcal{T}_{4^J}(4^j a) \in [0, 1/3] \quad \text{if } \theta_j = 0 \quad \text{and} \quad \mathcal{T}_{4^J}(4^j a) \in [2/3, 1] \quad \text{if } \theta_j = 1, \quad (3.6)$$

where  $\mathcal{T}_{4^J} : [0, 2 \times 4^J] \rightarrow [0, 1]$  is a “sawtooth” function with  $4^J$  “teeth” defined just above.

*Proof.* Fix  $j \in \{1, 2, \dots, J\}$ , we have

$$4^j a = 4^j \sum_{i=1}^J \theta_i 4^{-i} = \underbrace{\sum_{i=1}^{j-1} \theta_i 4^{j-i}}_{= 4k \text{ for some } k \in \mathbb{N} \text{ with } k \leq \frac{4^{j-1}}{3}} + \overbrace{\theta_j}^{0 \text{ or } 1} + \underbrace{\sum_{i=j+1}^J \theta_i 4^{j-i}}_{\in [0, \frac{1}{3}]}. \quad (3.7)$$

Clearly,

$$\sum_{i=1}^{j-1} \theta_i 4^{j-i} \in \{4k : k \in \mathbb{N}, k \leq 4^{j-1}/3\} \quad \text{and} \quad 0 \leq \sum_{i=j+1}^J \theta_i 4^{j-i} \leq \sum_{i=j+1}^J 4^{j-i} \leq 1/3.$$

If  $\theta_j = 0$ , then Equation (3.7) implies

$$4^j a \in [4k, 4k + 1/3] \quad \text{for some } k \in \mathbb{N} \text{ with } k \leq 4^{j-1}/3 \leq 4^j - 1,$$

which implies  $\mathcal{T}_{4^J}(4^j a) \in [0, 1/3]$ .

Similarly, if  $\theta_j = 1$ , then Equation (3.7) implies

$$4^j a \in [4k + 1, 4k + 1 + 1/3] \quad \text{for some } k \in \mathbb{N} \text{ with } k \leq 4^{j-1}/3 \leq 4^j - 1,$$

which implies  $\mathcal{T}_{4^J}(4^j a) \in [2/3, 1]$ . So we finish the proof.  $\square$

It is worth mentioning that the “sawtooth” function  $\mathcal{T}_m$  can be replaced by other functions that also have a key property “the function values near even integers are much larger than the ones near odd integers”.

With Lemma 3.2 in hand, we are ready to prove Proposition 3.1.

*Proof of Proposition 3.1.* By Lemma 3.2, we have

$$\mathcal{T}_{4^J}(4^j a) \in [0, 1/3] \quad \text{if } \theta_j = 0 \quad \text{and} \quad \mathcal{T}_{4^J}(4^j a) \in [2/3, 1] \quad \text{if } \theta_j = 1.$$

Define  $g(x) := 3\sigma(x - 1/3) - 3\sigma(x - 2/3)$  for any  $x \in \mathbb{R}$ , where  $\sigma$  is ReLU, i.e.,  $\sigma(x) = \max\{0, x\}$ . See an illustration of  $g$  in Figure 7. Clearly,

$$g(x) = 0 \text{ if } x \leq 1/3 \quad \text{and} \quad g(x) = 1 \text{ if } x \geq 2/3.$$

Hence,  $\phi := g \circ \mathcal{T}_{4^J}$  is the desired function. Obviously,  $\phi(x) \in [0, 1]$  for any  $x \in \mathbb{R}$ . To verify Equation (3.1), we fix  $j \in \{1, 2, \dots, J\}$ . If  $\theta_j = 0$ , then  $\mathcal{T}_{4^J}(4^j a) \in [0, 1/3]$ , implying  $\phi(4^j a) = g \circ \mathcal{T}_{4^J}(4^j a) = 0 = \theta_j$ . If  $\theta_j = 1$ , then  $\mathcal{T}_{4^J}(4^j a) \in [2/3, 1]$ , implying  $\phi(4^j a) = g \circ \mathcal{T}_{4^J}(4^j a) = 1 = \theta_j$ .

It remains to show  $\phi = g \circ \mathcal{T}_{4^J}$  can be realized by a ReLU network with the expected width and depth. Clearly,  $\mathcal{T}_{4^J}$  is a continuous piecewise linear function, which means

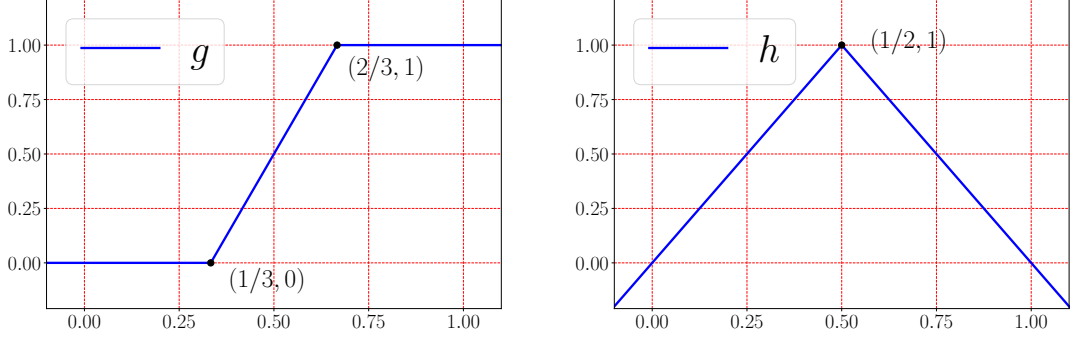


Figure 7: Illustrations of  $g$  and  $h$  on  $[0, 1]$ .

it can be implemented by a one-hidden-layer ReLU network. To make our construction more efficient, we introduce another method to implement  $\phi$ .

Define  $h(x) := 1 - 2\sigma(x - 1/2) - 2\sigma(1/2 - x)$  for any  $x \in [0, 1]$ . See an illustration of  $h$  in Figure 7. Then, it is easy to verify that  $\mathcal{T}_1(2x) = h(x)$  for any  $x \in [0, 1]$  and  $h$  can be implemented by a one-hidden-layer ReLU network with width 2. For any  $J \in \mathbb{N}^+$ , it is easy to verify that

$$\mathcal{T}_{4^J}(2^{2J+1}x) = \underbrace{h \circ h \circ \dots \circ h}_{2J+1 \text{ times}}(x), \quad \text{for any } x \in [0, 1].$$

So,  $\mathcal{T}_{4^J}$  limited on  $[0, 2 \times 4^J] = [0, 2^{2J+1}]$  can be realized by a ReLU network with width 2 and depth  $2J + 1$ . It follows that  $\phi = g \circ \mathcal{T}_{4^J}$  can be implemented by a ReLU network with width 2 and depth  $2J + 2$ , which means we finish the proof.  $\square$

## 4 Conclusion

This paper studies the approximation error of ReLU networks in terms of the number of intrinsic parameters. Theorem 1.1 implies that, for any Hölder continuous function on  $[0, 1]^d$ , a ReLU network with  $n + 2$  intrinsic parameters can approximate  $f$  with an error  $5\lambda d^{\alpha/2} 2^{-\alpha n}$  measured in the  $L^p$ -norm for  $p \in [1, \infty)$ , where  $\alpha \in (0, 1]$  and  $\lambda > 0$  are the Hölder order and constant, respectively. Moreover, such a result can be generalized from the  $L^p$ -norm to the  $L^\infty$ -norm with a price of adding  $\mathcal{O}(n)$  intrinsic parameters, as shown in Theorem 1.2. Finally, we show in Theorem 1.3 that three intrinsic parameters are enough to achieve an arbitrary error in the case of Hölder continuous functions, though this result requires high precision to encode these three parameters on computers. Remark that this paper only focuses on the approximation error characterized by the number of intrinsic parameters, the study of the optimization error and generalization error will be left as future work.

## Acknowledgments

Z. Shen is supported by Tan Chin Tuan Centennial Professorship. H. Yang was partially supported by the US National Science Foundation under award DMS-1945029.



S. Zhang is supported by a Postdoctoral Fellowship under NUS ENDOWMENT FUND (EXP WBS) (01 651).

## References

- [1] M. ANTHONY AND P. L. BARTLETT, *Neural Network Learning: Theoretical Foundations*, Cambridge University Press, New York, NY, USA, 1st ed., 2009.
- [2] C. BAO, Q. LI, Z. SHEN, C. TAI, L. WU, AND X. XIANG, *Approximation analysis of convolutional neural networks*, Semantic Scholar e-Preprint, (2019), p. Corpus ID: 204762668.
- [3] A. R. BARRON, *Universal approximation bounds for superpositions of a sigmoidal function*, IEEE Transactions on Information Theory, 39 (1993), pp. 930–945.
- [4] A. R. BARRON AND J. M. KLUSOWSKI, *Approximation and estimation for high-dimensional deep learning networks*, 2018.
- [5] P. BARTLETT, V. MAIOROV, AND R. MEIR, *Almost linear VC-dimension bounds for piecewise polynomial networks*, Neural Computation, 10 (1998), pp. 217–3.
- [6] J. BAXTER, *Theoretical models of learning to learn*, in Learning to Learn, Springer, 1998.
- [7] 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.
- [8] H. BÖLCSKEI, P. GROHS, G. KUTYNIOK, AND P. PETERSEN, *Optimal approximation with sparsely connected deep neural networks*, SIAM Journal on Mathematics of Data Science, 1 (2019), pp. 8–45.
- [9] S. BOZINOVSKI, *Reminder of the first paper on transfer learning in neural networks, 1976*, Informatica, 44 (2020), pp. 291–302.
- [10] S. BOZINOVSKI AND A. FULGOSI, *The influence of pattern similarity and transfer learning upon the training of a base perceptron B2. (original in croatian)*, in Proceedings of the Symposium Informatica 3-121-5, Bled, Slovenia, 1976.
- [11] L. CHEN AND C. WU, *A note on the expressive power of deep rectified linear unit networks in high-dimensional spaces*, Mathematical Methods in the Applied Sciences, 42 (2019), pp. 3400–3404.
- [12] M. CHEN, H. JIANG, W. LIAO, AND T. ZHAO, *Efficient approximation of deep relu networks for functions on low dimensional manifolds*, in Advances in Neural Information Processing Systems, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, eds., vol. 32, Curran Associates, Inc., 2019.

- [13] C. K. CHUI, S.-B. LIN, AND D.-X. ZHOU, *Construction of neural networks for realization of localized deep learning*, *Frontiers in Applied Mathematics and Statistics*, 4 (2018), p. 14.
- [14] G. CYBENKO, *Approximation by superpositions of a sigmoidal function*, *MCSS*, 2 (1989), pp. 303–314.
- [15] W. E, C. MA, AND L. WU, *A priori estimates of the population risk for two-layer neural networks*, *Communications in Mathematical Sciences*, 17 (2019), pp. 1407–1425.
- [16] R. GRIBONVAL, G. KUTYNIOK, M. NIELSEN, AND F. VOIGTLAENDER, *Approximation spaces of deep neural networks*, *arXiv e-prints*, (2019), p. arXiv:1905.01208.
- [17] I. GÜHRING, G. KUTYNIOK, AND P. PETERSEN, *Error bounds for approximations with deep ReLU neural networks in  $W^{s,p}$  norms*, *arXiv e-prints*, (2019), p. arXiv:1902.07896.
- [18] 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.
- [19] K. HORNIK, *Approximation capabilities of multilayer feedforward networks*, *Neural Networks*, 4 (1991), pp. 251–257.
- [20] K. HORNIK, M. STINCHCOMBE, AND H. WHITE, *Multilayer feedforward networks are universal approximators*, *Neural Networks*, 2 (1989), pp. 359–366.
- [21] Y. JIAO, Y. LAI, X. LU, F. WANG, J. ZHIJIAN YANG, AND Y. YANG, *Deep neural networks with ReLU-Sine-Exponential activations break curse of dimensionality on hölder class*, *arXiv e-prints*, (2021), p. arXiv:2103.00542.
- [22] T. G. KARIMPANAL AND R. BOUFFANAIS, *Self-organizing maps for storage and transfer of knowledge in reinforcement learning*, *Adaptive Behavior*, 27 (2019), pp. 111–126.
- [23] M. J. KEARNS AND R. E. SCHAPIRE, *Efficient distribution-free learning of probabilistic concepts*, *J. Comput. Syst. Sci.*, 48 (1994), pp. 464–497.
- [24] Q. LI, T. LIN, AND Z. SHEN, *Deep learning via dynamical systems: An approximation perspective*, *arXiv e-prints*, (2019), p. arXiv:1912.10382.
- [25] J. LU, Z. SHEN, H. YANG, AND S. ZHANG, *Deep network approximation for smooth functions*, *SIAM Journal on Mathematical Analysis*, 53 (2021), pp. 5465–5506.
- [26] L. MIHALKOVA, T. HUYNH, AND R. J. MOONEY, *Mapping and revising markov logic networks for transfer learning*, in *Proceedings of the 22nd Conference on Artificial Intelligence (AAAI-07)*, vol. 1, Vancouver, Canada, 2007, pp. 608–614.

- [27] T. MIKOLOV, A. JOULIN, AND M. BARONI, *A roadmap towards machine intelligence*, in Computational Linguistics and Intelligent Text Processing, A. Gelbukh, ed., Cham, 2018, Springer International Publishing, pp. 29–61.
- [28] H. MONTANELLI AND H. YANG, *Error bounds for deep ReLU networks using the Kolmogorov-Arnold superposition theorem*, Neural Networks, 129 (2020), pp. 1–6.
- [29] H. MONTANELLI, H. YANG, AND Q. DU, *Deep ReLU networks overcome the curse of dimensionality for bandlimited functions*, Journal of Computational Mathematics, (2020).
- [30] 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.
- [31] R. NAKADA AND M. IMAIZUMI, *Adaptive approximation and generalization of deep neural network with intrinsic dimensionality*, Journal of Machine Learning Research, 21 (2020), pp. 1–38.
- [32] A. NICULESCU-MIZIL AND R. CARUANA, *Inductive transfer for bayesian network structure learning*, in Proceedings of the Eleventh International Conference on Artificial Intelligence and Statistics, M. Meila and X. Shen, eds., vol. 2 of Proceedings of Machine Learning Research, San Juan, Puerto Rico, 21–24 Mar 2007, PMLR, pp. 339–346.
- [33] P. PETERSEN AND F. VOIGTLAENDER, *Optimal approximation of piecewise smooth functions using deep ReLU neural networks*, Neural Networks, 108 (2018), pp. 296–330.
- [34] L. Y. PRATT, *Discriminability-based transfer between neural networks*, in Advances in Neural Information Processing Systems, S. Hanson, J. Cowan, and C. Giles, eds., vol. 5, Morgan-Kaufmann, 1993.
- [35] A. A. RUSU, M. VECERÍK, T. ROTHÖRL, N. HEESS, R. PASCANU, AND R. HADSELL, *Sim-to-real robot learning from pixels with progressive nets*, in 1st Annual Conference on Robot Learning, CoRL 2017, Mountain View, California, USA, November 13–15, 2017, Proceedings, vol. 78 of Proceedings of Machine Learning Research, PMLR, 2017, pp. 262–270.
- [36] 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.
- [37] Z. SHEN, H. YANG, AND S. ZHANG, *Nonlinear approximation via compositions*, Neural Networks, 119 (2019), pp. 74–84.
- [38] ———, *Deep network approximation characterized by number of neurons*, Communications in Computational Physics, 28 (2020), pp. 1768–1811.

- [39] —, *Deep network approximation: Achieving arbitrary accuracy with fixed number of neurons*, arXiv e-prints, (2021), p. arXiv:2107.02397.
- [40] —, *Deep network with approximation error being reciprocal of width to power of square root of depth*, *Neural Computation*, 33 (2021), pp. 1005–1036.
- [41] —, *Neural network approximation: Three hidden layers are enough*, *Neural Networks*, 141 (2021), pp. 160–173.
- [42] —, *Optimal approximation rate of ReLU networks in terms of width and depth*, *Journal de Mathématiques Pures et Appliquées*, (to appear).
- [43] T. SUZUKI, *Adaptivity of deep ReLU network for learning in Besov and mixed smooth Besov spaces: optimal rate and curse of dimensionality*, in *International Conference on Learning Representations*, 2019.
- [44] D. YAROTSKY, *Error bounds for approximations with deep ReLU networks*, *Neural Networks*, 94 (2017), pp. 103–114.
- [45] —, *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.
- [46] D. YAROTSKY, *Elementary superexpressive activations*, in *Proceedings of the 38th International Conference on Machine Learning*, M. Meila and T. Zhang, eds., vol. 139 of *Proceedings of Machine Learning Research*, PMLR, 18–24 Jul 2021, pp. 11932–11940.
- [47] D. YAROTSKY AND A. ZHEVNERCHUK, *The phase diagram of approximation rates for deep neural networks*, in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, eds., vol. 33, Curran Associates, Inc., 2020, pp. 13005–13015.
- [48] 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>.
- [49] D.-X. ZHOU, *Universality of deep convolutional neural networks*, *Applied and Computational Harmonic Analysis*, 48 (2020), pp. 787–794.