

# COMPUTATION COMPLEXITY OF DEEP RELU NEURAL NETWORKS IN HIGH-DIMENSIONAL APPROXIMATION

DINH DŨNG<sup>1</sup>, VAN KIEN NGUYEN<sup>2,\*</sup>, MAI XUAN THAO<sup>3</sup>

<sup>1</sup>*Information Technology Institute, Vietnam National University, Hanoi  
144 Xuan Thuy, Cau Giay, Hanoi, Vietnam*

<sup>2</sup>*Faculty of Basic Sciences, University of Transport and Communications,  
No.3 Cau Giay Street, Lang Thuong Ward, Dong Da District, Hanoi, Vietnam*

<sup>3</sup>*Department of Natural Sciences, Hong Duc University,  
565 Quang Trung, Thanh Hoa, Vietnam*



**Abstract.** The purpose of the present paper is to study the computation complexity of deep ReLU neural networks for approximation of functions in Hölder-Nikol'skii spaces of mixed smoothness  $H_{\infty}^{\alpha}(\mathbb{I}^d)$  on the unit cube  $\mathbb{I}^d := [0, 1]^d$ . In this context, for any function  $f \in H_{\infty}^{\alpha}(\mathbb{I}^d)$ , we explicitly construct nonadaptive and adaptive deep ReLU neural networks having an output that approximates  $f$  with a prescribed accuracy  $\varepsilon$ , and prove dimension-dependent bounds for the computation complexity of the approximation, characterized by the size and depth of this deep ReLU neural network, explicitly in  $d$  and  $\varepsilon$ . Our results show the advantage of the adaptive method of approximation by deep ReLU neural networks over nonadaptive one.

**Keywords.** High-dimensional approximation; Deep ReLU neural networks; Nonadaptive approximation; Adaptive approximation; Hölder-Nikol'skii space of mixed smoothness.

## 1. INTRODUCTION

Neural networks and their applications have been of great interest for almost 80 years, started from the foundational works of McCulloch and Pitts [25], Hebb [19] and of Rosenblatt [31]. Recently, deep neural networks have been successfully applied to a wide variety of Machine Learning problems [22, 23, 38]. A typical deep neural network model is based on a hierarchy of composition of linear functions and a given nonlinear activation function. Among all the activation functions, Rectified Linear Unit (ReLU) activation  $\max(0, x)$  is the most commonly used in deep neural networks. Compared to traditional nonlinear activation functions (tanh and sigmoid), the ReLU function provides the same benefits but accelerates the training speed of deep neural networks since its gradient is very simple (either 0 or 1

---

Dedicated to Professor Phan Dinh Dieu on the occasion of his 85th birth anniversary.

\*Corresponding author.

*E-mail addresses:* dinhzung@gmail.com (D. Dũng); kiennv@utc.edu.vn (V.K. Nguyen); maixuanthao@hdu.edu.vn (M.X. Thao)

depending on the sign of input). Deep ReLU neural networks do not need to take additional time for computing error terms during training phase. Deep ReLU neural networks has found many applications in different fields such as computer vision [10, 18] and speech recognition [24, 35].

It is observed that there is a closed relation between approximating by deep ReLU networks and B-spline interpolation and quasi-interpolation representation, in particular by piecewise linear functions. There have been numerous papers addressing the expressive power of deep ReLU neural networks in approximating different sets of functions such as analytic functions [8, 26], differentiable functions [29, 39], oscillatory functions [14], functions in isotropic Sobolev or Besov spaces [1, 7, 11, 16, 40], functions with dominating mixed smoothness [27, 34] or in approximating solutions to partial differential equations [9, 28, 33], to mention just a few. The main advantage of deep ReLU neural networks in approximating functions is that they can output compositions of functions cheaply and consequently improve the convergence rate of approximation error, see [7, 8, 39]. We refer the reader to recent surveys [14, 30] for results in deep ReLU neural network approximation theory.

Function spaces having mixed smoothness appear naturally in many models of real world problem in mathematical physics, finance and other fields. For instance, in a recent work on regularity properties of solutions of the electronic Schrödinger equation, Yserentant [41] has shown that the eigenfunctions of the electronic Schrödinger operator have a certain mixed smoothness. Triebel [37, Chapter 6] has indicated a relation between Faber bases and sampling recovery in the context of spaces with mixed smoothness and solutions of Navier-Stokes equations. In particular, when initial data belongs to spaces with mixed smoothness, Navier-Stokes equations admit a unique solution. In mathematical finance, many problems are expressed as the expectation of some payoff function depending on quantities, such as stock prices, which are solutions of stochastic equations governed by Brownian motions. The payoff function normally has kinks and jumps and belongs to a very high dimensional space. To approximate the expected value one can apply preliminary integration method with respect to a single well chosen variable to obtain a function of  $d - 1$  variables which belongs to appropriate mixed Sobolev spaces in which Quasi-Monte Carlo can be applied efficiently, see [13] and references therein. For a survey on various aspects of high-dimensional approximation of functions having a mixed smoothness we refer the reader to the book [5].

In approximation theory, modern problems driven by a lot of applications in Information Technology, Mathematical Finance, Chemistry, Quantum Mechanics, Meteorology, and, in particular, in Uncertainty Quantification and Deep Machine Learning are being formulated in very high dimensions. Many times, numerical methods for such problems may demand computational cost increasing exponentially in dimension when the accuracy increases and as a consequence the method becomes intractable when the dimension of input data is large. This phenomenon is called “curse of dimensionality”. Hence, the problem of estimating dimension-dependent error in high-dimensional approximation problems arises naturally. Hyperbolic crosses and sparse grids promise to rid the “curse of dimensionality” in some problems when high-dimensional data belongs to certain classes of functions having mixed smoothness. Approximation methods and sampling algorithms for functions having mixed smoothness constructed on hyperbolic crosses and sparse grids give a surprising effect since hyperbolic crosses and sparse grids have the number of elements much less than those of standard domains and grids but give the same approximation error. This essentially reduces

the computational cost, and therefore makes the problem tractable.

In the recent paper [3], we have studied the approximation by deep ReLU neural networks, of functions from the Hölder-Zygmund space of mixed smoothness defined on the unit cube  $\mathbb{I}^d$  when the dimension  $d$  may be very large. The approximation error is measured in the norm of the isotropic Sobolev space. For any function  $f$  from this space, we explicitly constructed a deep ReLU neural network having an output that approximates  $f$  with a prescribed accuracy  $\varepsilon$ , and proved tight dimension-dependent estimates of the computation complexity of this approximation, characterized as the size and depth of this deep ReLU neural network, explicitly in  $d$  and  $\varepsilon$ .

The purpose of the present paper is to study the computation complexity of deep ReLU neural networks for approximation of Hölder-Nikol'skii functions having mixed smoothness on the unit cube  $\mathbb{I}^d := [0, 1]^d$ . Let us introduce the space  $H_\infty^\alpha(\mathbb{I}^d)$  of our interest. For univariate functions  $f$  on  $\mathbb{I} := [0, 1]$ , the difference operator  $\Delta_h$  is defined by

$$\Delta_h f(x) := f(x + h) - f(x),$$

for all  $x$  and  $h \geq 0$  such that  $x, x + h \in \mathbb{I}$ . If  $u$  is a subset of  $\{1, \dots, d\}$ , for multivariate functions  $f$  on  $\mathbb{I}^d$  the mixed difference operator  $\Delta_{\mathbf{h},u}$  is defined by

$$\Delta_{\mathbf{h},u} := \prod_{i \in u} \Delta_{h_i}, \quad \Delta_{\mathbf{h},\emptyset} = \text{Id},$$

for all  $\mathbf{x} = (x_1, \dots, x_d)$  and  $\mathbf{h} = (h_1, \dots, h_d)$  such that  $\mathbf{x}, \mathbf{x} + \mathbf{h} \in \mathbb{I}^d$ . Here the univariate operator  $\Delta_{h_i}$  is applied to the univariate function  $f$  by considering  $f$  as a function of variable  $x_i$  with the other variables held fixed. If  $0 < \alpha \leq 1$ , we introduce the semi-norm  $|f|_{H_\infty^\alpha(u)}$  for functions  $f \in C(\mathbb{I}^d)$  by

$$|f|_{H_\infty^\alpha(u)} := \sup_{\mathbf{h} > 0} \prod_{i \in u} h_i^{-\alpha} \|\Delta_{\mathbf{h},u}(f)\|_{C(\mathbb{I}^d(\mathbf{h},u))}$$

(in particular,  $|f|_{H_\infty^\alpha(\emptyset)} = \|f\|_{C(\mathbb{I}^d)}$ ), where  $\mathbb{I}^d(\mathbf{h}, u) := \{\mathbf{x} \in \mathbb{I}^d : x_i + h_i \in \mathbb{I}, i \in u\}$ . The Hölder-Nikol'skii space  $H_\infty^\alpha(\mathbb{I}^d)$  of mixed smoothness  $\alpha$  then is defined as the set of functions  $f \in C(\mathbb{I}^d)$  for which the norm

$$\|f\|_{H_\infty^\alpha(\mathbb{I}^d)} := \max_{u \subset \{1, \dots, d\}} |f|_{H_\infty^\alpha(u)}$$

is finite. From the definition we have that  $H_\infty^\alpha(\mathbb{I}^d) \subset C(\mathbb{I}^d)$ . Denote by  $\mathring{C}(\mathbb{I}^d)$  the set of all functions  $f \in C(\mathbb{I}^d)$  vanishing on the boundary  $\partial\mathbb{I}^d$  of  $\mathbb{I}^d$ , i.e., the set of all functions  $f \in C(\mathbb{I}^d)$  such that  $f(x) = 0$  if  $x_j = 0$  or  $x_j = 1$  for some index  $j \in \{1, \dots, d\}$ . Denote by  $\mathring{U}_\infty^{\alpha,d}$  the set of all functions  $f$  in the intersection  $\mathring{H}_\infty^\alpha(\mathbb{I}^d) := H_\infty^\alpha(\mathbb{I}^d) \cap \mathring{C}(\mathbb{I}^d)$  such that  $\|f\|_{H_\infty^\alpha(\mathbb{I}^d)} \leq 1$ .

As a continuation of [3], the present paper investigates nonadaptive and adaptive high-dimensional approximation by deep ReLU neural networks for functions from the classes  $\mathring{U}_\infty^{\alpha,d}$ . The approximation error is measured in the norm of  $L_\infty(\mathbb{I}^d)$ . In this context, we pay attention on the computation complexity of the deep ReLU networks, characterized by the size and depth of this deep ReLU neural network, explicitly in  $d$  and tolerance  $\varepsilon$ . A key tool for explicit construction of approximation methods by deep ReLU networks for functions in  $H_\infty^\alpha(\mathbb{I}^d)$  is truncations of tensorized Faber series.

The main contribution of the present paper is as follows.

Based on the decomposition of continuous functions by tensorized Faber series, for any  $f \in \mathring{U}_\infty^{\alpha,d}$  we explicitly construct a deep ReLU neural network  $\Phi_\varepsilon(f)$  having the output that

approximates  $f$  in the  $L_\infty(\mathbb{I}^d)$ -norm with a prescribed accuracy  $\varepsilon$  and having computation complexity expressing the dimension-dependent size

$$W(\Phi_\varepsilon(f)) \leq C_\alpha d \left( \frac{K_1^d}{(d-1)!} \right)^{\frac{1}{\alpha}+1} \varepsilon^{-\frac{1}{\alpha}} \log(2\varepsilon^{-1})^{(d-1)(\frac{1}{\alpha}+1)+1}, \tag{1}$$

where  $K_1 = B^{1/(\alpha+1)}4\alpha^{-1}$  with  $B = (2^\alpha - 1)^{-1}$ . The idea in proving the above result is to use truncation of Faber series  $R_n(f)$  as an intermediate approximation. Precisely, we first approximate function  $f \in \mathring{U}_\infty^{\alpha,d}$  by  $R_n(f)$  and then approximate  $R_n(f)$  by a deep ReLU network.

The advantage of this method is that the deep ReLU neural networks are easily constructed and they have the same architecture for all functions in  $\mathring{U}_\infty^{\alpha,d}$ , i.e., it is nonadaptive. However, since this method uses  $R_n(f)$  as an intermediate approximation, a disadvantage of it is that the computation complexity of deep ReLU networks is not better than that when approximating functions  $f \in \mathring{U}_\infty^{\alpha,d}$  by the linear method  $R_n(f)$ .

To overcome this disadvantage we develop a technique used in [40] and [7] for the univariate case. By this, we first represent the difference  $f - R_n(f)$  in a special form and then approximate terms in this representation by deep ReLU networks. We emphasize that extension of technique in [40] and [7] to multivariate case and mixed smoothness is non-trivial task since one needs to construct a set of finite cardinality to approximate functions in  $\mathring{U}_\infty^{\alpha,d}$ . For any  $f \in \mathring{U}_\infty^{\alpha,d}$  we explicitly construct a deep ReLU neural network  $\Phi_\varepsilon(f)$  of adaptive architecture having the output that approximates  $f$  in the  $L_\infty(\mathbb{I}^d)$ -norm with a prescribed accuracy  $\varepsilon$  and having a size estimated by

$$W(\Phi_\varepsilon(f)) \leq C_\alpha d^2 \left( \frac{K_2^d}{(d-1)!} \right)^{\frac{2}{\alpha}+2} \varepsilon^{-\frac{1}{\alpha}} (\log(2\varepsilon^{-1}) \log \log(2\varepsilon^{-1}))^{(1+\frac{1}{\alpha})(d-1)}, \tag{2}$$

where  $K_2 = 4(2^{\alpha+3}B)^{\frac{1}{2\alpha+2}}(\alpha^{-1} \log(2\alpha^{-1}))^{1/2}$ . Comparing (1) and (2) we find the later estimation improves  $\log(2\varepsilon^{-1})$ . Notice that the terms in right-hand side of both (1) and (2) which depend on dimension  $d$  only decay as fast as super exponential in  $d$ .

The outline of this paper is as follows. In Section 2, we introduce necessary definitions and elementary facts on deep ReLU neural networks. Section 3 is devoted to recall a decomposition of continuous functions on the unit cube  $\mathbb{I}^d$  by Faber system and approximation of functions  $f \in \mathring{U}_\infty^{\alpha,d}$  by truncations of Faber series  $R_n(f)$  as well as by sets of finite cardinality. In Section 4, we explicitly construct nonadaptive deep ReLU neural networks that approximate functions in  $\mathring{U}_\infty^{\alpha,d}$  and prove upper estimates for size and the depth required. Section 5 presents an improvement for approximation by adaptive deep ReLU neural networks of the results obtained in Section 4. In Section 6, we give an application of our results in numerical approximation of solutions to elliptic partial differential equations. Conclusions are given in Section 7.

**Notation.** As usual,  $\mathbb{N}$  is the natural numbers,  $\mathbb{Z}$  the integers,  $\mathbb{R}$  the real numbers and  $\mathbb{N}_0 := \{s \in \mathbb{Z} : s \geq 0\}$ ;  $\mathbb{N}_{-1} = \mathbb{N}_0 \cup \{-1\}$ . The letter  $d$  is reserved for the underlying dimension of  $\mathbb{R}^d$ ,  $\mathbb{N}^d$ , etc. If  $x \in \mathbb{R}$ ,  $[x]$  is defined to be the largest integer no larger than  $x$ . Vectorial quantities are denoted by boldface letters and  $x_i$  denotes the  $i$ th coordinate of  $\mathbf{x} \in \mathbb{R}^d$ , i.e.,  $\mathbf{x} := (x_1, \dots, x_d)$ . For  $\mathbf{x} \in \mathbb{R}^d$ , we denote  $|\mathbf{x}|_p := (|x_1|^p + \dots + |x_d|^p)^{1/p}$  if  $0 < p < \infty$ ,  $|\mathbf{x}|_\infty = \max\{|x_1|, \dots, |x_d|\}$ ,  $\text{supp}(\mathbf{x}) = \{j : x_j \neq 0\}$  and  $2^{\mathbf{x}} := (2^{x_1}, \dots, 2^{x_d})$ . If  $\mathbf{k}, \mathbf{s} \in \mathbb{N}_0^d$ , we denote  $2^{-\mathbf{k}}\mathbf{s} := (2^{-k_1}s_1, \dots, 2^{-k_d}s_d)$ . We use the abbreviation:  $L_\infty := L_\infty(\mathbb{I}^d)$

and  $\|\cdot\|_\infty := \|\cdot\|_{L_\infty}$ . Universal constants or constants depending on parameters  $\alpha, d$  are denoted by  $C$  or  $C_{\alpha,d}$ , respectively. Values of constants  $C$  and  $C_{\alpha,d}$  in general, are not specified except the case when they are precisely given, and may be different in various places.  $|A|$  denotes the cardinality of the finite set  $A$ .

## 2. DEEP RELU NEURAL NETWORKS

In this section we introduce necessary definitions and elementary facts on deep ReLU neural networks. There is a wide variety of neural network architectures and each of them is adapted to specific tasks. We only consider feed-forward deep ReLU neural networks such that only connections between neighboring layers are allowed.

**Definition 1.** Let  $d, L \in \mathbb{N}$ ,  $L \geq 2$ ,  $N_0 = d$ , and  $N_1, \dots, N_L \in \mathbb{N}$ . Let  $\mathbf{W}^\ell = (w_{i,j}^\ell)$ ,  $\ell = 1, \dots, L$ , be  $N_\ell \times N_{\ell-1}$  matrix, and  $\mathbf{b}^\ell = (b_j^\ell) \in \mathbb{R}^{N_\ell}$ .

- A neural network  $\Phi$  with input dimension  $d$  and  $L$  layers is a sequence of matrix-vector tuples

$$\Phi = ((\mathbf{W}^1, \mathbf{b}^1), \dots, (\mathbf{W}^L, \mathbf{b}^L)).$$

We will use the following terminology.

- The number of layers  $L(\Phi) = L$  is the depth of  $\Phi$ ;
- $N_w(\Phi) = \max_{\ell=0, \dots, L} \{N_\ell\}$  is the width of  $\Phi$ ;  $\mathbf{N}(\Phi) = (N_0, N_1, \dots, N_L)$  the dimension of  $\Phi$ ;
- The real numbers  $w_{i,j}^\ell$  and  $b_j^\ell$  are edge and node weights of  $\Phi$ , respectively;
- The number of nonzero weights  $w_{i,j}^\ell$  and  $b_j^\ell$  is the size of  $\Phi$  and denoted by  $W(\Phi)$ ;
- When  $L(\Phi) \geq 3$ ,  $\Phi$  is called a deep neural network, and otherwise, a shallow neural network.
- A neural network architecture  $\mathbb{A}$  with input dimension  $d$  and  $L$  layers is a neural network

$$\mathbb{A} = ((\mathbf{W}^1, \mathbf{b}^1), \dots, (\mathbf{W}^L, \mathbf{b}^L)),$$

where elements of  $\mathbf{W}^\ell$  and  $\mathbf{b}^\ell$ ,  $\ell = 1, \dots, L$ , are in  $\{0, 1\}$ .

The above defined networks are sometimes called standard networks to distinguish with networks allowing for connections of neurons in non-neighboring layers. A deep neural network can be visualized in a graph. The graph associated with a deep neural network  $\Phi$  defined in Definition 1 consists of  $L + 1$  layers which are numbered from 0 to  $L$ . The  $\ell$ th layer has  $N_\ell$  nodes which are numbered from 1 to  $N_\ell$ . If  $w_{i,j}^\ell \neq 0$ , then there is an edge connecting the node  $j$  in the layer  $\ell - 1$  to the node  $i$  in the layer  $\ell$ . In Figure 1 we illustrate a deep neural network with input dimension 3 and 5 layers.

**Definition 2.** Given  $L \in \mathbb{N}$ ,  $L \geq 2$ , and a neural network architecture  $\mathbb{A} = ((\overline{\mathbf{W}}^1, \overline{\mathbf{b}}^1), \dots, (\overline{\mathbf{W}}^L, \overline{\mathbf{b}}^L))$ . We say that a neural network  $\Phi = ((\mathbf{W}^1, \mathbf{b}^1), \dots, (\mathbf{W}^L, \mathbf{b}^L))$  has architecture  $\mathbb{A}$  if

- $\mathbf{N}(\Phi) = \mathbf{N}(\mathbb{A})$ ;
- $\overline{w}_{i,j}^\ell = 0$  implies  $w_{i,j}^\ell = 0$ ,  $\overline{b}_i^\ell = 0$  implies  $b_i^\ell = 0$  for all  $i = 1, \dots, N_\ell$ ,  $j = 1, \dots, N_{\ell-1}$ , and  $\ell = 1, \dots, L$ . Here  $\overline{w}_{i,j}^\ell$  are entries of  $\overline{\mathbf{W}}^\ell$  and  $\overline{b}_i^\ell$  are elements of  $\overline{\mathbf{b}}^\ell$ ,  $\ell = 1, \dots, L$ .

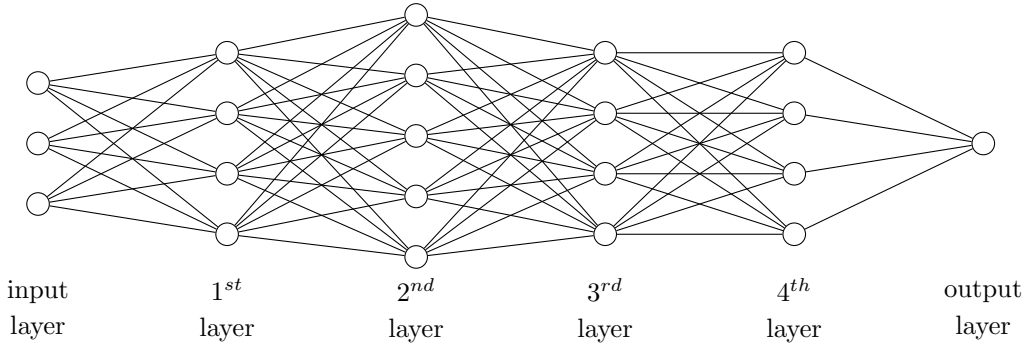


Figure 1: The graph associated to a deep neural network with input dimension 3 and 5 layers

For a given deep neural network  $\Phi = ((\mathbf{W}^1, \mathbf{b}^1), \dots, (\mathbf{W}^L, \mathbf{b}^L))$ , there exists a unique deep neural network architecture  $\mathbb{A} = ((\bar{\mathbf{W}}^1, \bar{\mathbf{b}}^1), \dots, (\bar{\mathbf{W}}^L, \bar{\mathbf{b}}^L))$  such that

- $\mathcal{N}(\Phi) = \mathcal{N}(\mathbb{A})$ ;
- $\bar{w}_{i,j}^\ell = 0 \iff w_{i,j}^\ell = 0, \bar{b}_i^\ell = 0 \iff b_i^\ell = 0$  for all  $i = 1, \dots, N_\ell, j = 1, \dots, N_{\ell-1}$ , and  $\ell = 1, \dots, L$ .

We call this architecture  $\mathbb{A}$  the minimal architecture of  $\Phi$  (this definition is proper in the sense that any architecture of  $\Phi$  is also an architecture of  $\mathbb{A}$ ).

A deep neural network is associated with an activation function which calculates output at each node. The choice of activation function depends on the problem under consideration. In this paper we focus our attention on ReLU activation function defined by  $\sigma(t) := \max\{t, 0\}, t \in \mathbb{R}$ . We will use the notation  $\sigma(\mathbf{x}) := (\sigma(x_1), \dots, \sigma(x_d))$  for  $\mathbf{x} \in \mathbb{R}^d$ .

**Definition 3.** A deep ReLU neural network with input dimension  $d$  and  $L$  layers is a neural network

$$\Phi = ((\mathbf{W}^1, \mathbf{b}^1), \dots, (\mathbf{W}^L, \mathbf{b}^L)),$$

in which the following computation scheme is implemented

$$\begin{aligned} \mathbf{z}^0 &:= \mathbf{x} \in \mathbb{R}^d, \\ \mathbf{z}^\ell &:= \sigma(\mathbf{W}^\ell \mathbf{z}^{\ell-1} + \mathbf{b}^\ell), \quad \ell = 1, \dots, L-1, \\ \mathbf{z}^L &:= \mathbf{W}^L \mathbf{z}^{L-1} + \mathbf{b}^L. \end{aligned}$$

We call  $\mathbf{z}^0$  the input and with an ambiguity denote  $\Phi(\mathbf{x}) := \mathbf{z}^L$  the output of  $\Phi$  and in some places we identify a deep ReLU neural network with its output.

Several deep ReLU neural networks can be combined to form a larger deep ReLU neural network whose output is a linear combination or composition of outputs of sub-networks. In the following, we introduce parallelization, concatenation and special construction.

**Lemma 1** (Parallelization). *Let  $N \in \mathbb{N}, \Omega \subset \mathbb{R}^d$  be a bounded set,  $\lambda_j \in \mathbb{R}, j = 1, \dots, N$ . Let  $\Phi_j, j = 1, \dots, N$  be deep ReLU neural networks with input dimension  $d$ . Then we can explicitly construct a deep ReLU neural network denoted by  $\Phi$  so that*

$$\Phi(\mathbf{x}) = \sum_{j=1}^N \lambda_j \Phi_j(\mathbf{x}), \quad \mathbf{x} \in \Omega,$$

with  $L(\Phi) = \max_{j=1, \dots, N} \{L(\Phi_j)\}$  and

$$W(\Phi) = \sum_{j=1}^N W(\Phi_j) + \sum_{j:L(\Phi_j) < L(\Phi)} (L(\Phi) - L(\Phi_j) + 2) \leq 3N \max_{j=1, \dots, N} W(\Phi_j).$$

The network  $\Phi$  is called the *Parallelization network* of  $\Phi_j$ ,  $j = 1, \dots, N$ .

A proof of Lemma 1 can be found in [3]. The last estimate in Lemma 1 is due to  $2 \leq L \leq \max_{j=1, \dots, N} W(\Phi_j)$ .

Another way to construct a ReLU network whose output is a linear combination of outputs of other ReLU networks is to use special networks. A special deep ReLU neural network with input dimension  $d$  can be defined as follows. In each hidden layer a special role is reserved for  $d$  first (top) nodes and the last (bottom) node. Concatenation of top  $d$  nodes and the bottom node in each layer to the corresponding nodes in the next layer form  $d + 1$  parallel channels. The nodes in these  $d + 1$  channel are free of activation. The top  $d$  parallel channels are called the source channels and just carry  $\mathbf{x} = (x_1, \dots, x_d)$  forward. The bottom channel is called collation channel. The nodes in the bottom channel are used to collect intermediate outputs by addition. This channel never feeds forward into subsequent calculation, it only accepts previous calculations. It has been shown in [3] that if  $\Phi$  is a special deep ReLU neural network with input dimension  $d$  depth  $L$  and  $\mathbf{x} \in \mathbb{I}^d$ , then there is a deep ReLU neural network  $\Phi'$  such that

$$L(\Phi') = L(\Phi) \tag{3}$$

and  $\Phi'(\mathbf{x}) = \Phi(\mathbf{x})$ . In view of the proof of [3, Lemma 4.2] we find only node weights in the collation channel of  $\Phi$  and  $\Phi'$  are different. Therefore we deduce

$$W(\Phi') \leq W(\Phi) + L(\Phi) \leq 2W(\Phi). \tag{4}$$

**Lemma 2** (Special Construction). *Let  $N \in \mathbb{N}$ ,  $\Omega \subset \mathbb{R}^d$  be a bounded set,  $\lambda_j \in \mathbb{R}$ ,  $j = 1, \dots, N$ . Let  $\Phi_j$ ,  $j = 1, \dots, N$  be deep ReLU neural networks with input dimension  $d$ . Then we can explicitly construct a deep special ReLU neural network denoted by  $\Phi$  so that*

$$\Phi(\mathbf{x}) = \sum_{j=1}^N \lambda_j \Phi_j(\mathbf{x}), \quad \mathbf{x} \in \Omega,$$

with  $L(\Phi) = \sum_{j=1}^N L(\Phi_j)$  and  $W(\Phi) \leq \sum_{j=1}^N W(\Phi_j) + (d + 1)L(\Phi)$ .

An illustration of a special network  $\Phi$  whose output is a linear combination of network  $\Phi_j$ ,  $j = 1, \dots, N$  is given in Figure 2.

**Lemma 3** (Concatenation). *Let  $\Phi_1$  and  $\Phi_2$  be two ReLU neural networks such that output layer of  $\Phi_1$  has the same dimension as input layer of  $\Phi_2$ . Then, we can explicitly construct a ReLU neural network  $\Phi$  such that  $\Phi(\mathbf{x}) = \Phi_2(\Phi_1(\mathbf{x}))$  for  $\mathbf{x} \in \mathbb{R}^d$ . Moreover we have  $L(\Phi) = L(\Phi_1) + L(\Phi_2)$  and  $W(\Phi) \leq 2W(\Phi_1) + 2W(\Phi_2)$ .*

A proof of the above lemma can be found in [29]. The network  $\Phi$  in this lemma is called the concatenation network of  $\Phi_1$  and  $\Phi_2$ .

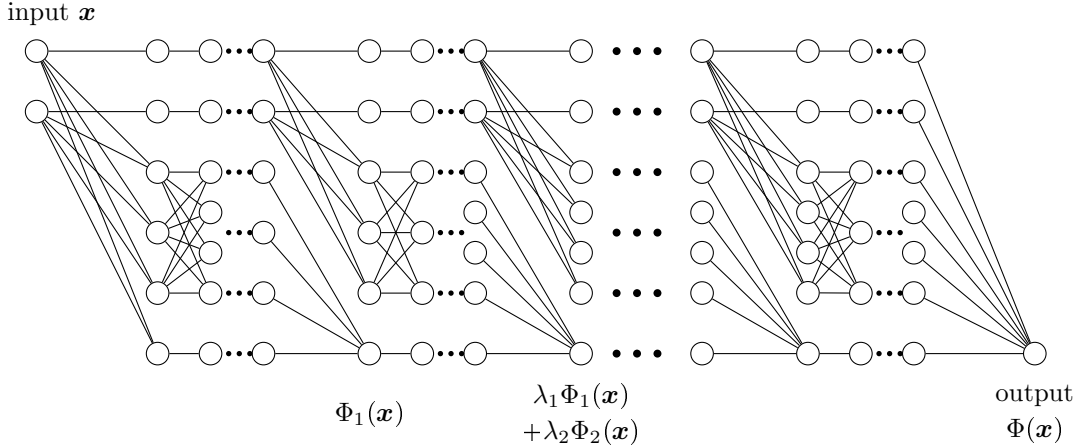


Figure 2: Illustration of a special deep ReLU neural network ( $d = 2$ )

### 3. APPROXIMATION BY SETS OF FINITE CARDINALITY

In this section we recall a decomposition of continuous functions on the unit cube  $\mathbb{I}^d$  by Faber series, interpolation approximation by truncated Faber series and by set of finite cardinality. They are a key tool for explicit construction of approximation methods by deep ReLU networks for functions in Hölder-Nikol'skii spaces of mixed smoothness.

Let  $\varphi(x) = (1 - |x - 1|)_+$ ,  $x \in \mathbb{R}$ , be the hat function (the piece-wise linear B-spline with knots at  $0, 1, 2$ ), where  $x_+ := \max(x, 0)$  for  $x \in \mathbb{R}$ . For  $k \in \mathbb{N}_{-1}$  we define the functions  $\varphi_{k,s}$  by

$$\varphi_{k,s}(x) := \varphi(2^{k+1}x - 2s), \quad k \geq 0, \quad s \in Z(k) := \{0, 1, \dots, 2^k - 1\},$$

and

$$\varphi_{-1,s}(x) := \varphi(x - s + 1), \quad s \in Z(-1) := \{0, 1\}.$$

For a univariate function  $f$  on  $\mathbb{I}$ ,  $k \in \mathbb{N}_{-1}$ , and  $s \in Z(k)$  we define

$$\lambda_{k,s}(f) := -\frac{1}{2} \Delta_{2^{-k-1}}^2 f(2^{-k}s), \quad k \geq 0, \quad \lambda_{-1,s}(f) := f(s),$$

where

$$\Delta_h^2 f(x) := f(x + 2h) - 2f(x + h) + f(x),$$

for all  $x$  and  $h \geq 0$  such that  $x, x + h \in \mathbb{I}$ . If  $m \in \mathbb{N}_0$  we put

$$R_m(f) := \sum_{k=0}^m q_k(f), \quad q_k(f) := \sum_{s \in Z(k)} \lambda_{k,s}(f) \varphi_{k,s}. \tag{5}$$

For  $k \in \mathbb{N}_0$ , we define the functions  $\varphi_{k,s}^* \in \mathring{C}(\mathbb{I})$  by

$$\varphi_{k,s}^*(x) := \varphi(2^{k+1}x - s + 1), \quad s \in Z_*(k) := \{1, \dots, 2^{k+1} - 1\}, \tag{6}$$

and for  $f \in \mathring{C}(\mathbb{I})$  one can check

$$R_m(f) = \sum_{s \in Z_*(m)} f(2^{-m-1}s) \varphi_{m,s}^*.$$



Hence  $R_m(f) \in \mathring{C}(\mathbb{I})$  interpolates  $f$  at the points  $2^{-m-1}\mathbf{s}$ ,  $\mathbf{s} \in Z_*(m)$ , that is,

$$R_m(f)(2^{-m-1}\mathbf{s}) = f(2^{-m-1}\mathbf{s}), \quad \mathbf{s} \in Z_*(m).$$

Put  $Z(\mathbf{k}) := \times_{j=1}^d Z(k_j)$ . For  $\mathbf{k} \in \mathbb{N}_{-1}^d$ ,  $\mathbf{s} \in Z(\mathbf{k})$ , we introduce the tensorized Faber basis by

$$\varphi_{\mathbf{k},\mathbf{s}}(\mathbf{x}) := \varphi_{k_1,s_1}(x_1) \cdot \dots \cdot \varphi_{k_d,s_d}(x_d), \quad \mathbf{x} \in \mathbb{I}^d. \tag{7}$$

We also define the linear functionals  $\lambda_{\mathbf{k},\mathbf{s}}$  for multivariate function  $f$  on  $\mathbb{I}^d$ ,  $\mathbf{k} \in \mathbb{N}_{-1}^d$ , and  $\mathbf{s} \in Z(\mathbf{k})$  by

$$\lambda_{\mathbf{k},\mathbf{s}}(f) := \prod_{i=1}^d \lambda_{k_i,s_i}(f),$$

where the univariate functional  $\lambda_{k_i,s_i}$  is applied to the univariate function  $f$  by considering  $f$  as a function of variable  $x_i$  with the other variables held fixed.

We have the following representation by Faber series for continuous functions on  $\mathbb{I}^d$  (see [2, Section 4] and [36, Theorem 3.10]).

**Lemma 4.** *The tensorized Faber system  $\{\varphi_{\mathbf{k},\mathbf{s}} : \mathbf{k} \in \mathbb{N}_{-1}^d, \mathbf{s} \in Z(\mathbf{k})\}$  is a basis in  $C(\mathbb{I}^d)$ . Moreover, every function  $f \in C(\mathbb{I}^d)$  can be represented by the Faber series*

$$f = \sum_{\mathbf{k} \in \mathbb{N}_{-1}^d} q_{\mathbf{k}}(f), \quad q_{\mathbf{k}}(f) := \sum_{\mathbf{s} \in Z(\mathbf{k})} \lambda_{\mathbf{k},\mathbf{s}}(f) \varphi_{\mathbf{k},\mathbf{s}}$$

converging in the norm of  $C(\mathbb{I}^d)$ .

When  $f \in \mathring{U}_{\infty}^{\alpha,d}$ ,  $\lambda_{\mathbf{k},\mathbf{s}}(f) = 0$  if  $k_j = -1$  for some  $j \in \{1, \dots, d\}$ , hence we can write

$$f = \sum_{\mathbf{k} \in \mathbb{N}_0^d} q_{\mathbf{k}}(f)$$

with unconditional convergence in  $C(\mathbb{I}^d)$ , see [36, Theorem 3.13]. In this case it holds the following estimate

$$\begin{aligned} |\lambda_{\mathbf{k},\mathbf{s}}(f)| &= 2^{-d} \left| \prod_{i=1}^d \Delta_{2^{-k_i-1}}^2 f(2^{-\mathbf{k}}\mathbf{s}) \right| \\ &= 2^{-d} \left| \prod_{i=1}^d \left[ \Delta_{2^{-k_i-1}} f(2^{-\mathbf{k}}\mathbf{s} + 2^{-k_i-1}\mathbf{e}^i) - \Delta_{2^{-k_i-1}} f(2^{-\mathbf{k}}\mathbf{s}) \right] \right| \leq 2^{-\alpha d} 2^{-\alpha|\mathbf{k}|_1}, \end{aligned} \tag{8}$$

for  $\mathbf{k} \in \mathbb{N}_0^d$ ,  $\mathbf{s} \in Z(\mathbf{k})$ . Here  $\{\mathbf{e}^i\}_{i=1,\dots,d}$  is the standard basis of  $\mathbb{R}^d$ .

For  $f \in \mathring{C}(\mathbb{I}^d)$ , we define the operator  $R_m$  by

$$R_m(f) := \sum_{|\mathbf{k}|_1 \leq m} q_{\mathbf{k}}(f) = \sum_{|\mathbf{k}|_1 \leq m} \sum_{\mathbf{s} \in Z(\mathbf{k})} \lambda_{\mathbf{k},\mathbf{s}}(f) \varphi_{\mathbf{k},\mathbf{s}}.$$

The truncated Faber series  $R_m(f) \in \mathring{C}(\mathbb{I}^d)$  completely determined by values of  $f$  at the points  $2^{-\mathbf{k}-1}\mathbf{s}$ , for  $(\mathbf{k}, \mathbf{s}) \in G^d(m)$ , where

$$G^d(m) := \{(\mathbf{k}, \mathbf{s}) : |\mathbf{k}|_1 \leq m, \mathbf{s} \in Z_*(\mathbf{k})\},$$

$Z_*(\mathbf{k}) := \prod_{j=1}^d Z_*(k_j)$  and  $\mathbf{1} = (1, \dots, 1) \in \mathbb{N}^d$ . Moreover,  $R_m(f)$  interpolates  $f$  at the points  $2^{-\mathbf{k}-1}\mathbf{s}$ , for  $(\mathbf{k}, \mathbf{s}) \in G^d(m)$ , i.e.,

$$R_m(f)(2^{-\mathbf{k}-1}\mathbf{s}) = f(2^{-\mathbf{k}-1}\mathbf{s}), \quad (\mathbf{k}, \mathbf{s}) \in G^d(m).$$

The following lemma gives a  $d$ -dependent estimate of the approximation error by  $R_m(f)$  of  $f \in \mathring{U}_\infty^{\alpha,d}$ , see [4].

**Lemma 5.** *Let  $d \geq 2$ ,  $m \in \mathbb{N}$ , and  $0 < \alpha \leq 1$ . Then we have*

$$\sup_{f \in \mathring{U}_\infty^{\alpha,d}} \|f - R_m(f)\|_\infty \leq 2^{-\alpha} B^d 2^{-\alpha m} \binom{m+d}{d-1}, \quad B = (2^\alpha - 1)^{-1}.$$

We make use the abbreviations:  $\mathbf{x}_j := (x_1, \dots, x_j) \in \mathbb{R}^j$ ;  $\bar{\mathbf{x}}_j := (x_{j+1}, \dots, x_d) \in \mathbb{R}^{d-j}$  with the convention  $\mathbf{x}_0 := 0$  for  $\mathbf{x} \in \mathbb{R}^d$  and  $j = 0, 1, \dots, d-1$ . When  $j = 1$  we denote  $x_1$  instead of  $\mathbf{x}_1$ .

For  $f \in \mathring{U}_\infty^{\alpha,1}$  we explicitly construct the function  $S_f \in \mathring{C}(\mathbb{I})$  by

$$S_f := \sum_{s \in Z_*(m)} 2^{-\alpha(m+1)} l_s(f) \varphi_{m,s}^*, \tag{9}$$

where we put  $l_0(f) = 0$  and assign the values  $S_f(2^{-m-1}s) = 2^{-\alpha(m+1)} l_s(f)$  from left to right closest to  $f(2^{-m-1}s)$  for  $s = 1, \dots, 2^{m+1} - 1$ . If there are two possible choices for  $l_s(f)$  we choose  $l_s(f)$  that is closest to the already determined  $l_{s-1}(f)$ . We define

$$\mathcal{S}^\alpha(m) := \{S_f : f \in \mathring{U}_\infty^{\alpha,1}\}. \tag{10}$$

It has been proved that the set  $\mathcal{S}^\alpha(m)$  is finite and it holds the estimate  $|\mathcal{S}^\alpha(m)| \leq 3^{2^{m+1}}$ , see [4]. Moreover, by Lemma 5 and [4, Lemma 2.3] for  $f \in \mathring{U}_\infty^{\alpha,1}$  and  $m \in \mathbb{N}_0$  we have

$$\|f - S_f\|_\infty \leq \|f - R_m(f)\|_\infty + \|R_m(f) - S_f\|_\infty \leq 2^{-(m+1)\alpha - \frac{1}{2}} + \frac{2^{-(m+1)\alpha}}{2^\alpha - 1}. \tag{11}$$

In case of high dimensions we have the following.

**Lemma 6.** *Let  $m > 1$ ,  $d \geq 2$  and  $0 < \alpha \leq 1$ . For  $f \in \mathring{U}_\infty^{\alpha,d}$ , let the function  $S_m(f)$  be defined by*

$$S_m(f)(\mathbf{x}) := \sum_{|\bar{\mathbf{k}}_1| \leq m} 2^{-\alpha(|\bar{\mathbf{k}}_1|+d-1)} \sum_{\bar{\mathbf{s}}_1 \in Z(\bar{\mathbf{k}}_1)} \varphi_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1}(\bar{\mathbf{x}}_1) S_{K_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1}}(f)(x_1), \tag{12}$$

where  $S_{K_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1}}(f) \in \mathcal{S}^\alpha(m - |\bar{\mathbf{k}}_1|_1)$  is as in (9) for the function  $K_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1}(f)$ . Then it holds the inequality

$$\|f - S_m(f)\|_\infty \leq B^d 2^{-\alpha m} \binom{m+d}{d-1}.$$

Moreover, for the set

$$\mathcal{S}^{\alpha,d}(m) := \{S_m(f) : f \in \mathring{U}_\infty^{\alpha,d}\},$$

we have  $N_d(m) := |\mathcal{S}^{\alpha,d}(m)| \leq 3^{2^{m+1} \binom{m+d-1}{d-1}}$ .

For a proof of the above lemma we refer the reader to [4].

#### 4. DEEP RELU NEURAL NETWORK APPROXIMATION A NONADAPTIVE METHOD

In this section, we explicitly construct a nonadaptive deep ReLU neural network having an output that approximates every function  $f \in \dot{U}_\infty^{\alpha,d}$  in the  $L_\infty(\mathbb{I}^d)$ -norm with a prescribed accuracy  $\varepsilon$  and prove dimension-dependent error estimates of its size and depth. Nonadaptivity means that its architecture is the same for all  $f \in \dot{U}_\infty^{\alpha,d}$ . Our technique is first to approximate  $f$  by its truncation of Faber series  $R_n(f)$  and then approximate  $R_n(f)$  by a deep ReLU network. Since the case  $d = 1$  was already considered (see, e.g., [1, 7, 11]), we study the high dimension case when  $d \geq 2$ . Our main result in this section is read as follows.

**Theorem 7.** *Let  $d \in \mathbb{N}$ ,  $d \geq 2$  and  $\alpha \in (0, 1]$ . Then there is  $\varepsilon_0 = \varepsilon_0(d, \alpha) \in (0, 1]$  such that for every  $\varepsilon \in (0, \varepsilon_0)$  we can explicitly construct a deep neural network architecture  $\mathbb{A}_\varepsilon$  with the following property. For every  $f \in \dot{U}_\infty^{\alpha,d}$ , we can explicitly construct a deep ReLU neural network  $\Phi_\varepsilon(f)$  having the architecture  $\mathbb{A}_\varepsilon$  such that*

$$\|f - \Phi_\varepsilon(f)\|_\infty \leq \varepsilon.$$

Moreover, we have

$$W(\mathbb{A}_\varepsilon) \leq C_\alpha d \left( \frac{K_1^d}{(d-1)!} \right)^{\frac{1}{\alpha}+1} \varepsilon^{-\frac{1}{\alpha}} \log(2\varepsilon^{-1})^{(d-1)(\frac{1}{\alpha}+1)+1} \tag{13}$$

and

$$L(\mathbb{A}_\varepsilon) \leq C \log d \log(2\varepsilon^{-1}),$$

where  $K_1 = B^{1/(\alpha+1)}4\alpha^{-1}$  with  $B$  given in Lemma 5 and  $C_\alpha$  depends only on  $\alpha$ .

To prepare for proving Theorem 7 we recall results of approximating the product  $\prod_{j=1}^d x_j$  and  $\varphi_{\mathbf{k},\mathbf{s}}$  by deep ReLU neural networks, see [33] and [3].

**Lemma 8.** *For every  $\delta \in (0, 1)$ ,  $d \in \mathbb{N}$ ,  $d \geq 2$ , we can explicitly construct a deep ReLU neural network  $\Phi_P$  so that*

$$\sup_{\mathbf{x} \in [-1,1]^d} \left| \prod_{j=1}^d x_j - \Phi_P(\mathbf{x}) \right| \leq \delta.$$

Furthermore, if  $x_j = 0$  for some  $j \in \{1, \dots, d\}$  then  $\Phi_P(\mathbf{x}) = 0$  and there exists a constant  $C > 0$  independent of  $\delta$  and  $d$  such that

$$W(\Phi_P) \leq Cd \log(d\delta^{-1}) \quad \text{and} \quad L(\Phi_P) \leq C \log d \log(d\delta^{-1}).$$

**Lemma 9.** *For every dimension  $d \geq 2$ ,  $\delta \in (0, 1)$  and for the  $d$ -variate hat functions  $\varphi_{\mathbf{k},\mathbf{s}}$ ,  $\mathbf{k} \in \mathbb{N}_0^d$ ,  $\mathbf{s} \in Z(\mathbf{k})$ , defined as in (7), we can explicitly construct a deep neural network  $\Phi_\delta(\varphi_{\mathbf{k},\mathbf{s}})$  so that*

$$\|\varphi_{\mathbf{k},\mathbf{s}} - \Phi_\delta(\varphi_{\mathbf{k},\mathbf{s}})\|_\infty \leq \delta$$

and

$$W(\Phi_\delta(\varphi_{\mathbf{k},\mathbf{s}})) \leq Cd \log(d\delta^{-1}) \quad \text{and} \quad L(\Phi_\delta(\varphi_{\mathbf{k},\mathbf{s}})) \leq C \log d \log(d\delta^{-1}). \tag{14}$$

Moreover,  $\text{supp } \Phi_\delta(\varphi_{\mathbf{k},\mathbf{s}}) \subset \text{supp } \varphi_{\mathbf{k},\mathbf{s}}$ .

The above result allows us to construct a deep ReLU network  $\Phi_\varepsilon(R_n(f))$  to approximate  $R_n(f)$ .

**Lemma 10.** *Let  $d \in \mathbb{N}$ ,  $d \geq 2$ ,  $n \in \mathbb{N}$ ,  $\alpha \in (0, 1]$  and  $\varepsilon \in (0, 1)$ . Then for every  $f \in \dot{U}_\infty^{\alpha, d}$  we can explicitly construct a deep ReLU network  $\Phi_\varepsilon(R_n(f))$  of the same architecture  $\mathbb{A}_\varepsilon$  so that*

$$\|R_n(f) - \Phi_\varepsilon(R_n(f))\|_\infty \leq \varepsilon. \quad (15)$$

Moreover, we have

$$W(\Phi_\varepsilon(R_n(f))) \leq Cd2^n \binom{n+d-1}{d-1} \log(dB^d\varepsilon^{-1}) \quad (16)$$

and

$$L(\Phi_\varepsilon(R_n(f))) \leq C \log d \log(dB^d\varepsilon^{-1}). \quad (17)$$

The estimates (16) and (17) also hold for  $W(\mathbb{A}_\varepsilon)$  and  $L(\mathbb{A}_\varepsilon)$  respectively.

*Proof.* For every pair  $\mathbf{k}, \mathbf{s}$  with  $|\mathbf{k}|_1 \leq n$  and  $\mathbf{s} \in Z(\mathbf{k})$ , by applying Lemma 9 with  $\delta := B^{-d}\varepsilon$ , we explicitly construct a deep ReLU neural network  $\Phi_\delta(\varphi_{\mathbf{k}, \mathbf{s}})$  so that  $\text{supp } \Phi_\delta(\varphi_{\mathbf{k}, \mathbf{s}}) \subset \text{supp } \varphi_{\mathbf{k}, \mathbf{s}}$ ,

$$\|\varphi_{\mathbf{k}, \mathbf{s}} - \Phi_\delta(\varphi_{\mathbf{k}, \mathbf{s}})\|_\infty \leq B^{-d}\varepsilon, \quad (18)$$

and it holds the estimates (14) for  $W(\Phi_\delta(\varphi_{\mathbf{k}, \mathbf{s}}))$  and  $L(\Phi_\delta(\varphi_{\mathbf{k}, \mathbf{s}}))$ . We approximate  $R_n(f)$  by the output

$$\Phi_\varepsilon(R_n(f))(\mathbf{x}) = \sum_{|\mathbf{k}|_1 \leq n} \sum_{\mathbf{s} \in Z(\mathbf{k})} \lambda_{\mathbf{k}, \mathbf{s}}(f) \Phi_\delta(\varphi_{\mathbf{k}, \mathbf{s}})(\mathbf{x})$$

of the network  $\Phi_\varepsilon(R_n(f))$  which is a parallelization of the networks  $\{\Phi_\delta(\varphi_{\mathbf{k}, \mathbf{s}})\}_{|\mathbf{k}|_1 \leq n, \mathbf{s} \in Z(\mathbf{k})}$ . Notice that the interiors of  $\text{supp } \Phi_\delta(\varphi_{\mathbf{k}, \mathbf{s}})$  and  $\text{supp } \Phi_\delta(\varphi_{\mathbf{k}, \mathbf{s}'})$  have empty intersection if  $\mathbf{s} \neq \mathbf{s}'$ . Moreover, for every  $\mathbf{x} \in \mathbb{I}^d$ , there is an  $\mathbf{s} \in Z(\mathbf{k})$  such that  $\mathbf{x} \in \text{supp } \varphi_{\mathbf{k}, \mathbf{s}}$ , and hence, by using (8) and (18) we get the estimates

$$\begin{aligned} |R_n(f)(\mathbf{x}) - \Phi_\varepsilon(R_n(f))(\mathbf{x})| &= \sum_{|\mathbf{k}|_1 \leq n} |\lambda_{\mathbf{k}, \mathbf{s}}(f)(\varphi_{\mathbf{k}, \mathbf{s}}(\mathbf{x}) - \Phi_\delta(\varphi_{\mathbf{k}, \mathbf{s}})(\mathbf{x}))| \\ &\leq 2^{-\alpha d} \sum_{|\mathbf{k}|_1 \leq n} 2^{-\alpha|\mathbf{k}|_1} \varepsilon B^{-d} \\ &\leq \varepsilon(1 - 2^{-\alpha})^d \sum_{j=0}^n 2^{-\alpha j} \binom{j+d-1}{d-1}. \end{aligned}$$

From

$$\sum_{j=0}^{\infty} \binom{j+m}{m} t^j \leq (1-t)^{-m-1}, \quad t \in (0, 1), \quad (19)$$

see [6, Lemma 2.2], we obtain (15).

By using Lemma 1 and the estimates (14), the size and the depth of  $\Phi_\varepsilon(R_n(f))$  can be estimated as

$$\begin{aligned} W(\Phi_\varepsilon(R_n(f))) &\leq C |\{(\mathbf{k}, \mathbf{s}) : |\mathbf{k}|_1 \leq n, \mathbf{s} \in Z(\mathbf{k})\}| \max_{|\mathbf{k}|_1 \leq n, \mathbf{s} \in Z(\mathbf{k})} W(\Phi_\delta(\varphi_{\mathbf{k}, \mathbf{s}})) \\ &= C \sum_{\ell=0}^n 2^\ell \binom{\ell+d-1}{d-1} d \log(dB^d\varepsilon^{-1}) \\ &\leq Cd2^n \binom{n+d-1}{d-1} \log(dB^d\varepsilon^{-1}), \end{aligned}$$

and

$$L(\Phi_\varepsilon(R_n(f))) \leq \max_{|\mathbf{k}|_1 \leq n, \mathbf{s} \in Z(\mathbf{k})} L(\Phi_\delta(\varphi_{\mathbf{k}, \mathbf{s}})) \leq C \log d \log(dB^d \varepsilon^{-1}).$$

The proof is completed by noticing that  $\Phi_\varepsilon(R_n(f))$  has the architecture  $\mathbb{A}_\varepsilon$  (independent of  $f$ ) which is defined as the minimal architecture of the deep ReLU neural network  $\Phi_\varepsilon$  obtained by parallelization of the networks  $\{\Phi_\delta(\varphi_{\mathbf{k}, \mathbf{s}})\}_{|\mathbf{k}|_1 \leq n, \mathbf{s} \in Z(\mathbf{k})}$  with the output

$$\Phi_\varepsilon(\mathbf{x}) = \sum_{|\mathbf{k}|_1 \leq n} \sum_{\mathbf{s} \in Z(\mathbf{k})} \Phi_\delta(\varphi_{\mathbf{k}, \mathbf{s}})(\mathbf{x}), \quad \mathbf{x} \in \mathbb{I}^d.$$

Hence, the estimates (16) and (17) also hold for  $W(\mathbb{A}_\varepsilon)$  and  $L(\mathbb{A}_\varepsilon)$  respectively.  $\blacksquare$

We are ready to prove Theorem 7.

*Proof.* Denote  $n_0$  the natural point from which the function  $h(n) = 2^{-\alpha} B^d 2^{-\alpha n} \binom{n+d}{d-1}$  is decreasing and  $h(n-1) \leq 2^{-\alpha n/2}$  for all  $n > n_0$ . We put  $\varepsilon_0 = \min\{h(n_0), h(d)\}$ . For  $\varepsilon \in (0, \varepsilon_0)$  we define  $n > \max\{n_0, d\}$  by

$$2^{-\alpha} B^d 2^{-\alpha n} \binom{n+d}{d-1} \leq \frac{\varepsilon}{2} < 2^{-\alpha} B^d 2^{-\alpha(n-1)} \binom{n-1+d}{d-1}. \quad (20)$$

With  $\varepsilon' = \varepsilon/2$  in Lemma 10 and  $\Phi_\varepsilon(f) = \Phi_{\varepsilon'}(R_n(f))$  we have

$$\|f - \Phi_\varepsilon(f)\|_\infty \leq \|f - R_n(f)\|_\infty + \|R_n - \Phi_{\varepsilon'}(R_n(f))\|_\infty \leq 2^{-\alpha} B^d 2^{-\alpha n} \binom{n+d}{d-1} + \frac{\varepsilon}{2} \leq \varepsilon.$$

We define  $\mathbb{A}_\varepsilon$  as the minimal architecture of the deep ReLU neural network  $\Phi_\varepsilon$  obtained by parallelization of the networks  $\{\Phi_\delta(\varphi_{\mathbf{k}, \mathbf{s}})\}_{|\mathbf{k}|_1 \leq n, \mathbf{s} \in Z(\mathbf{k})}$  with the output

$$\Phi_\varepsilon(\mathbf{x}) = \sum_{|\mathbf{k}|_1 \leq n} \sum_{\mathbf{s} \in Z(\mathbf{k})} \Phi_\delta(\varphi_{\mathbf{k}, \mathbf{s}})(\mathbf{x}), \quad \mathbf{x} \in \mathbb{I}^d.$$

Then  $\Phi_\varepsilon(f)$  has the architecture for all  $f \in \dot{U}_\infty^{\alpha, d}$ . From Lemma 10 we have

$$W(\mathbb{A}_\varepsilon) \leq Cd2^n \log(2d\varepsilon^{-1}B^d) \binom{n+d-1}{d-1}.$$

From the choice of  $n$  we have

$$2d\varepsilon^{-1}B^d \leq d2^\alpha 2^{\alpha n} \binom{n+d}{d-1}^{-1} \leq 2^{\alpha n} 2d \binom{n+d}{d-1}^{-1} \leq 2^{\alpha n}.$$

By this and (20) we get

$$\begin{aligned} W(\mathbb{A}_\varepsilon) &\leq Cd \left( 2\varepsilon^{-1} B^d \binom{n+d-1}{d-1} \right)^{1/\alpha} \alpha n \binom{n+d-1}{d-1} \\ &\leq Cd (2\varepsilon^{-1} B^d)^{1/\alpha} n \binom{n+d-1}{d-1}^{\frac{1}{\alpha}+1} \\ &\leq Cd (\varepsilon^{-1} B^d)^{1/\alpha} n \left( \frac{(2n)^{d-1}}{(d-1)!} \right)^{\frac{1}{\alpha}+1}. \end{aligned}$$

Now  $h(n-1) \leq 2^{-\alpha n/2}$  and (20) lead to  $\frac{\varepsilon}{2} \leq 2^{-\alpha n/2}$  which implies  $n \leq \frac{2}{\alpha} \log(2\varepsilon^{-1})$ . Therefore we get

$$\begin{aligned} W(\mathbb{A}_\varepsilon) &\leq Cd(\varepsilon^{-1}B^d)^{1/\alpha} \log(2\varepsilon^{-1}) \left( \frac{(4\alpha^{-1} \log(2\varepsilon^{-1}))^{d-1}}{(d-1)!} \right)^{\frac{1}{\alpha}+1} \\ &= Cd(B^d)^{1/\alpha} \left( \frac{(4\alpha^{-1})^{d-1}}{(d-1)!} \right)^{\frac{1}{\alpha}+1} \varepsilon^{-\frac{1}{\alpha}} \log(2\varepsilon^{-1})^{(d-1)(\frac{1}{\alpha}+1)+1}, \end{aligned}$$

and (13) follows. We also have

$$L(\mathbb{A}_\varepsilon) \leq C \log d \log(d2\varepsilon^{-1}B^d) \leq C\alpha n \log d \leq C \log d \log(2\varepsilon^{-1}).$$

■

## 5. DEEP RELU NEURAL NETWORK APPROXIMATION - AN ADAPTIVE METHOD

In this section, we explicitly construct an adaptive method of approximation with accuracy  $\varepsilon > 0$  by deep ReLU neural networks of functions  $f \in \mathring{U}_\infty^{\alpha,d}$ . This method reduces the computation complexity expressing as the size and depth of the approximating deep ReLU networks comparing with the computation complexity of the nonadaptive method given in Theorem 7. As mentioned the univariate case was already considered in [7] ( $0 < \alpha < 1$ ) and [40] ( $\alpha = 1$ ), we focus our attention on multivariate case when  $d \geq 2$ . The main result of this section is read as follows.

**Theorem 11.** *Let  $d \in \mathbb{N}$ ,  $d \geq 2$ ,  $\alpha \in (0, 1]$ . Then there is  $\varepsilon_0 = \varepsilon_0(d, \alpha) \in (0, 1/2]$  such that for every  $\varepsilon \in (0, \varepsilon_0)$  and for every  $f \in \mathring{U}_\infty^{\alpha,d}$  we can explicitly construct an adaptive deep ReLU neural network  $\Phi_\varepsilon(f)$  so that*

$$\|f - \Phi_\varepsilon(f)\|_\infty \leq \varepsilon.$$

Moreover, we have

$$W(\Phi_\varepsilon(f)) \leq C_\alpha d^2 \left( \frac{K_2^d}{(d-1)!} \right)^{\frac{2}{\alpha}+2} \varepsilon^{-\frac{1}{\alpha}} (\log(2\varepsilon^{-1}) \log \log(2\varepsilon^{-1}))^{(1+\frac{1}{\alpha})(d-1)} \quad (21)$$

and

$$L(\Phi_\varepsilon(f)) \leq C'_\alpha \varepsilon^{-\frac{1}{d\alpha}} (\log(2\varepsilon^{-1}))^{\frac{d-1-\alpha}{d\alpha}} (\log \log(2\varepsilon^{-1}))^{\frac{(\alpha+1)(d-1)}{d\alpha}}, \quad (22)$$

where

$$K_2 := 4(2^{\alpha+3}B)^{\frac{1}{2\alpha+2}} (\alpha^{-1} \log(2\alpha^{-1}))^{1/2}$$

and positive constants  $C_\alpha, C'_\alpha$  depend on  $\alpha$  only.

Let us explain the idea of the proof. Let  $f \in \mathring{U}_\infty^{\alpha,d}$  and  $\varepsilon \in (0, \varepsilon_0)$  ( $\varepsilon_0$  will be specified latter) be given. Using the writing

$$f = R_n(f) + (f - R_n(f)),$$

we explicitly construct deep ReLU neural networks to approximate with accuracy  $\varepsilon/2$  the terms  $R_n(f)$  and  $f - R_n(f)$  and evaluate the dimension-dependent computation complexity

separately, and then take their sum to get an approximation with accuracy  $\varepsilon$  to  $f$  and its dimension-dependent computation complexity. For approximation of the first term  $R_n(f)$ , we take the deep ReLU neural network  $\Phi_{\varepsilon/2}(R_n(f))$  which has been constructed in Lemma 10.

Thus, our main task is to explicitly construct a desired deep ReLU neural network  $\Phi_{\varepsilon/2}(f - R_n(f))$  for approximation of the second term  $f - R_n(f)$ . Our strategy is to represent the difference  $f - R_n(f)$  in a special form and then approximate terms in this representation by deep ReLU networks. To this end, we need some auxiliary preparation.

For univariate functions  $f \in \mathring{C}(\mathbb{I})$ , let the operator  $T_k$ ,  $k \in \mathbb{N}_0$ , be defined by

$$T_k(f) := f - R_{k-1}(f)$$

with the operator  $R_k$  defined as in (5) and the convention  $R_{-1} := 0$ . From this definition we have  $T_0$  is the identity operator. Notice that for  $f \in \mathring{U}_\infty^{\alpha,1}$ , it holds the inequality  $\|T_k(f)\|_{H_\infty^\alpha(\mathbb{I})} \leq 2$ .

For a multivariate function  $f \in \mathring{C}(\mathbb{I}^d)$ , the tensor product operator  $T_{\mathbf{k}}$ ,  $\mathbf{k} = (k_1, \dots, k_d) \in \mathbb{N}_0^d$ , is defined by

$$T_{\mathbf{k}}(f) := \prod_{j=1}^d T_{k_j}(f),$$

where the univariate operator  $T_{k_j}$  is applied to the univariate function  $f$  by considering  $f$  as a function of variable  $x_j$  with the other variables held fixed.

For  $n \in \mathbb{N}$ , it has been shown in [4] that  $f - R_n(f)$  can be represented in the following special form

$$f - R_n(f) = \sum_{j=0}^{d-1} \sum_{|\mathbf{k}_j|_1 \leq n} F_{\mathbf{k}_j}, \tag{23}$$

where  $F_{\mathbf{k}_0} := T_{(n+1)\mathbf{e}^1}$  and

$$F_{\mathbf{k}_j} := T_{(n+1-|\mathbf{k}_j|_1)\mathbf{e}^{j+1}}(q_{\mathbf{k}_j}(f)), \quad j = 1, \dots, d-1,$$

or equivalently,

$$F_{\mathbf{k}_j} = \prod_{i=1}^j (T_{(k_i-1)\mathbf{e}^j} - T_{k_i\mathbf{e}^j}) T_{(n+1-|\mathbf{k}_j|_1)\mathbf{e}^{j+1}}(f), \quad j = 1, \dots, d-1. \tag{24}$$

We shall explicitly construct deep ReLU neural networks  $\Phi_{\varepsilon'}(F_{\mathbf{k}_j})$  to approximate each term  $F_{\mathbf{k}_j}$  in the sum in (23). Due to (24) this is reduced to construct deep ReLU networks that approximate  $T_{\mathbf{k}}(f)$ ,  $\mathbf{k} \in \mathbb{N}_0^d$ . Put

$$I_{\mathbf{k},\mathbf{s}} := \bigotimes_{j=1}^d I_{k_j,s_j} = \bigotimes_{j=1}^d [2^{-k_j} s_j, 2^{-k_j}(s_j + 1)], \quad \mathbf{k} \in \mathbb{N}_0^d, \quad \mathbf{s} \in Z(\mathbf{k}),$$

and

$$T_{\mathbf{k},\mathbf{s}}(f)(\mathbf{x}) := 2^{\alpha|\mathbf{k}|-d} (T_{\mathbf{k}}(f)\chi_{I_{\mathbf{k},\mathbf{s}}})(2^{-\mathbf{k}}(\mathbf{x} + \mathbf{s})).$$

Since  $\text{supp}(T_{\mathbf{k}}(f)\chi_{I_{\mathbf{k},\mathbf{s}}}) \subset I_{\mathbf{k},\mathbf{s}}$  and  $\|T_{\mathbf{k}}(f)\chi_{I_{\mathbf{k},\mathbf{s}}}\|_{H_\infty^\alpha(\mathbb{I}^d)} \leq 2^d$ , we have that

$$\text{supp}(T_{\mathbf{k},\mathbf{s}}(f)) \subset \mathbb{I}^d, \quad T_{\mathbf{k},\mathbf{s}}(f) \in \mathring{U}_\infty^{\alpha,d}.$$

Take the function  $S_m(T_{\mathbf{k},s}(f))$  defined as in (12) for  $T_{\mathbf{k},s}(f) \in \dot{U}_\infty^{\alpha,d}$ . By Lemma 6 it holds the estimate

$$\|T_{\mathbf{k},s}(f) - S_m(T_{\mathbf{k},s}(f))\|_\infty \leq B^d 2^{-\alpha m} \binom{m+d}{d-1}.$$

Define

$$S_{\mathbf{k},m}(f)(\mathbf{x}) := 2^{-\alpha|\mathbf{k}|_1+d} \sum_{\mathbf{s} \in Z(\mathbf{k})} S_m(T_{\mathbf{k},s}(f))(2^{\mathbf{k}}\mathbf{x} - \mathbf{s}). \quad (25)$$

We then get

$$\begin{aligned} \|T_{\mathbf{k}}(f) - S_{\mathbf{k},m}(f)\|_\infty &= \left\| \sum_{\mathbf{s} \in Z(\mathbf{k})} \left[ T_{\mathbf{k}}(f) \chi_{I_{\mathbf{k},s}}(\cdot) - 2^{-\alpha|\mathbf{k}|_1+d} S_m(T_{\mathbf{k},s}(f))(2^{\mathbf{k}}\cdot - \mathbf{s}) \right] \right\|_\infty \\ &= 2^{-\alpha|\mathbf{k}|_1+d} \left\| \sum_{\mathbf{s} \in Z(\mathbf{k})} \left[ T_{\mathbf{k},s}(f) - S_m(T_{\mathbf{k},s}(f)) \right] (2^{\mathbf{k}}\cdot - \mathbf{s}) \right\|_\infty. \end{aligned}$$

Since support of  $T_{\mathbf{k},s}(f) - S_m(T_{\mathbf{k},s}(f))$  is contained in  $\mathbb{I}^d$ , we finally obtain

$$\|T_{\mathbf{k}}(f) - S_{\mathbf{k},m}(f)\|_\infty \leq (2B)^d (2^m 2^{|\mathbf{k}|_1})^{-\alpha} \binom{m+d}{d-1}. \quad (26)$$

Considering  $S_{\mathbf{k},m}(f)$  as an intermediate approximation of  $T_{\mathbf{k}}(f)$ , we shall construct deep ReLU networks approximating  $S_{\mathbf{k},m}(f)$ . Since  $S_{\mathbf{k},m}(f)$  is a sum of functions in  $\mathcal{S}^{\alpha,d}(m)$ , we shall construct a deep ReLU neural network  $\Phi_\varepsilon(S)$  for approximating  $S \in \mathcal{S}^{\alpha,d}(m)$  with accuracy  $\varepsilon$  and estimate its size.

**Lemma 12.** *Let  $d \in \mathbb{N}$ ,  $d \geq 2$ ,  $m \in \mathbb{N}$ ,  $\alpha \in (0, 1]$ , and  $\varepsilon \in (0, 1)$ . Then for every  $S \in \mathcal{S}^{\alpha,d}(m)$ , we can explicitly construct a deep ReLU neural network  $\Phi_\varepsilon(S)$  so that  $\text{supp } \Phi_\varepsilon(S) \subset \mathbb{I}^d$  and*

$$\|S - \Phi_\varepsilon(S)\|_\infty \leq \varepsilon. \quad (27)$$

Moreover, there is a positive constant  $C$  such that

$$W(\Phi_\varepsilon(S)) \leq Cd \log d 2^m \binom{m+d-1}{d-1} \log(dB^d \varepsilon^{-1}) \quad (28)$$

and

$$L(\Phi_\varepsilon(S)) \leq C 2^m \log d \log(dB^d \varepsilon^{-1}), \quad (29)$$

where  $B$  is given in Lemma 5.

*Proof.* By Lemma 6, for every function  $S \in \mathcal{S}^{\alpha,d}(m)$ , there is a function  $f \in \dot{U}_\infty^{\alpha,d}$  such that

$$S(\mathbf{x}) = S_m(f)(\mathbf{x}) = \sum_{|\bar{\mathbf{k}}_1|_1 \leq m} 2^{-\alpha(|\bar{\mathbf{k}}_1|_1+d-1)} \sum_{\bar{\mathbf{s}}_1 \in Z(\bar{\mathbf{k}}_1)} \varphi_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1}(\bar{\mathbf{x}}_1) S_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1}(x_1),$$

where  $S_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1} := S_{K_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1}}(f) \in \mathcal{S}^\alpha(m - |\bar{\mathbf{k}}_1|_1)$ . Since  $S_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1}$  is a piecewise linear continuous function, see (10) and (9), according to [7, Theorem 3.1] we can explicitly construct a deep ReLU neural network  $\Phi(S_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1})$  with one-dimensional input so that  $\Phi(S_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1})(x_1) = S_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1}(x_1)$ ,  $x_1 \in \mathbb{I}$ , and

$$W(\Phi(S_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1})) \leq C 2^{m-|\bar{\mathbf{k}}_1|_1}, \quad L(\Phi(S_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1})) \leq C 2^{m-|\bar{\mathbf{k}}_1|_1}. \quad (30)$$



Each univariate function  $\varphi_{k_j, s_j}$  in the tensor product  $\varphi_{\bar{k}_1, \bar{s}_1} = \otimes_{j=2}^d \varphi_{k_j, s_j}$  can be expressed as an output of a neural network  $\Phi(\varphi_{k_j, s_j})$  with one-dimensional input, deep 3 and 8 weights. Adding layers (with one node in each layer) putting forward  $x_j$  to each network  $\Phi(\varphi_{k_j, s_j})$  such that it has the length  $L(\Phi(S_{\bar{k}_1, \bar{s}_1}))$ . We still denote these new networks by  $\Phi(\varphi_{k_j, s_j})$ . Then we obtain

$$W(\Phi(\varphi_{k_j, s_j})) \leq C2^{m-|\bar{k}_1|_1}.$$

We approximate the  $d$ -univariate function  $\varphi_{\bar{k}_1, \bar{s}_1}(\bar{x}_1)S_{\bar{k}_1, \bar{s}_1}(x_1)$  by the output of the network  $\Phi_{\bar{k}_1, \bar{s}_1}$  with  $d$ -dimensional input which is explicitly constructed as a concatenation of the networks  $\Phi(S_{\bar{k}_1, \bar{s}_1})$ ,  $\Phi(\varphi_{k_j, s_j})$ ,  $j = 2, \dots, d$ , with product network  $\Phi_P$  in Lemma 8. With  $\delta = \varepsilon B^{1-d}$  in Lemma 8 we have

$$\|\varphi_{\bar{k}_1, \bar{s}_1}S_{\bar{k}_1, \bar{s}_1} - \Phi_{\bar{k}_1, \bar{s}_1}\|_\infty \leq \varepsilon B^{1-d}. \quad (31)$$

Since  $|\varphi_{\bar{k}_1, \bar{s}_1}(\bar{x}_1)| \leq 1$  for  $\bar{x}_1 \in \mathbb{I}^{d-1}$  and  $|S_{\bar{k}_1, \bar{s}_1}(x_1)| \leq 4$  for  $x_1 \in \mathbb{I}$  by (11), from Lemmata 3, 8 and (30) we derive that

$$\begin{aligned} W(\Phi_{\bar{k}_1, \bar{s}_1}) &\leq C \left( \sum_{j=2}^d W(\Phi(\varphi_{k_j, s_j})) + W(\Phi(S_{\bar{k}_1, \bar{s}_1})) + W(\Phi_P) \right) \\ &\leq Cd(2^{m-|\bar{k}_1|_1} + \log(dB^d\varepsilon^{-1})), \end{aligned} \quad (32)$$

and

$$L(\Phi_{\bar{k}_1, \bar{s}_1}) \leq L(\Phi(S_{\bar{k}_1, \bar{s}_1})) + L(\Phi_P) \leq C(2^{m-|\bar{k}_1|_1} + \log d \log(dB^d\varepsilon^{-1})). \quad (33)$$

Moreover  $\text{supp}(\Phi_{\bar{k}_1, \bar{s}_1}) \subset \text{supp}(\varphi_{\bar{k}_1, \bar{s}_1}S_{\bar{k}_1, \bar{s}_1})$  by Lemma 8.

Let the network  $\Phi_{\bar{k}_1}$  with output

$$\Phi_{\bar{k}_1}(\mathbf{x}) = \sum_{\bar{s}_1 \in Z(\bar{k}_1)} \Phi_{\bar{k}_1, \bar{s}_1}(\mathbf{x})$$

be explicitly constructed as a combination of the networks  $\{\Phi_{\bar{k}_1, \bar{s}_1}\}_{\bar{s}_1 \in Z(\bar{k}_1)}$  by the special construction. Then by Lemma 2, (32) and (33) we obtain that

$$\begin{aligned} L(\Phi_{\bar{k}_1}) &\leq \sum_{\bar{s}_1 \in Z(\bar{k}_1)} L(\Phi_{\bar{k}_1, \bar{s}_1}) \leq C2^{|\bar{k}_1|_1} (2^{m-|\bar{k}_1|_1} + \log d \log(dB^d\varepsilon^{-1})) \\ &\leq C2^m \log d \log(dB^d\varepsilon^{-1}) \end{aligned} \quad (34)$$

and

$$\begin{aligned} W(\Phi_{\bar{k}_1}) &\leq \sum_{\bar{s}_1 \in Z(\bar{k}_1)} W(\Phi_{\bar{k}_1, \bar{s}_1}) + (d+1)L(\Phi_{\bar{k}_1}) \\ &\leq Cd2^{|\bar{k}_1|_1} (2^{m-|\bar{k}_1|_1} + \log(dB^d\varepsilon^{-1})) + C(d \log d)2^m \log(dB^d\varepsilon^{-1}) \\ &\leq C(d \log d)2^m \log(dB^d\varepsilon^{-1}). \end{aligned} \quad (35)$$

Since  $\mathbf{x} \in \mathbb{I}^d$ , we can construct a standard network with the same output as  $\Phi_{\bar{k}_1}$  and the estimates (34) and (35) hold, see (3) and (4). We still denote this network by  $\Phi_{\bar{k}_1}$ . Now we define the network  $\Phi_\varepsilon(S)$  as a parallelization of the networks  $(\Phi_{\bar{k}_1})_{|\bar{k}_1|_1 \leq m}$  with output

$$\Phi_\varepsilon(S)(\mathbf{x}) = \sum_{|\bar{k}_1|_1 \leq m} 2^{-\alpha(|\bar{k}_1|_1 + d - 1)} \Phi_{\bar{k}_1}(\mathbf{x}).$$

Since  $\text{supp}(\Phi_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1}) \subset \text{supp}(\varphi_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1} S_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1})$  and for a given  $\bar{\mathbf{k}}_1$ ,  $\text{supp}(\varphi_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1} S_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1})$  and  $\text{supp}(\varphi_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}'_1} S_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}'_1})$  are disjoint if  $\bar{\mathbf{s}}'_1 \neq \bar{\mathbf{s}}_1$ , it holds by (31) and (19) that

$$\begin{aligned} \|S - \Phi_\varepsilon(S)\|_\infty &\leq \left\| \sum_{|\bar{\mathbf{k}}_1|_1 \leq m} \sum_{\bar{\mathbf{s}}_1 \in Z(\bar{\mathbf{k}}_1)} 2^{-\alpha(|\bar{\mathbf{k}}_1|_1 + d - 1)} |\varphi_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1} S_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1} - \Phi_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1}| \right\|_\infty \\ &= \sum_{|\bar{\mathbf{k}}_1|_1 \leq m} 2^{-\alpha(|\bar{\mathbf{k}}_1|_1 + d - 1)} \max_{\bar{\mathbf{s}}_1 \in Z(\bar{\mathbf{k}}_1)} \|\varphi_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1} S_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1} - \Phi_{\bar{\mathbf{k}}_1, \bar{\mathbf{s}}_1}\|_\infty \\ &\leq \sum_{|\bar{\mathbf{k}}_1|_1 \leq m} 2^{-\alpha(|\bar{\mathbf{k}}_1|_1 + d - 1)} \varepsilon B^{1-d} \\ &= \varepsilon (1 - 2^{-\alpha})^{d-1} \sum_{\ell=0}^m 2^{-\ell\alpha} \binom{\ell + d - 2}{d - 2} \leq \varepsilon. \end{aligned}$$

By Lemma 1 and (34), (35) we obtain

$$W(\Phi_\varepsilon(S)) \leq 3 \left| \{\bar{\mathbf{k}}_1 : |\bar{\mathbf{k}}_1|_1 \leq m\} \right| \max_{|\bar{\mathbf{k}}_1|_1 \leq m} W(\Phi_{\bar{\mathbf{k}}_1}) \leq C(d \log d) 2^m \binom{m + d - 1}{d - 1} \log(dB^d \varepsilon^{-1}),$$

and

$$L(\Phi_\varepsilon(S)) \leq \max_{|\bar{\mathbf{k}}_1|_1 \leq m} L(\Phi_{\bar{\mathbf{k}}_1}) \leq C(\log d) 2^m \log(dB^d \varepsilon^{-1}).$$

Finally, the inclusion  $\text{supp} \Phi_\varepsilon(S) \subset \mathbb{I}^d$  follows from Lemmata 8 and 9.  $\blacksquare$

The following result is a generalization of [7, Lemma 5.1] to  $d$ -dimensional case.

**Lemma 13.** *Let  $k \in \mathbb{N}$ ,  $\Lambda \subset Z(k)$  and  $j \in \{1, \dots, d\}$ . Let  $\Phi$  be a deep ReLU network with input dimension  $d$  such that  $\text{supp} \Phi \subset \mathbb{I}^d$ . Denote*

$$f(\mathbf{x}) := \sum_{s \in \Lambda} \Phi(x_1, \dots, 2^k x_j - s, \dots, x_d), \quad \mathbf{x} \in \mathbb{I}^d.$$

*Then we can explicitly construct a deep ReLU network  $\Phi_\Lambda$  with output  $f(\mathbf{x})$  and*

$$W(\Phi_\Lambda) \leq C(d|\Lambda| + W(\Phi)), \quad L(\Phi_\Lambda) \leq 5 + L(\Phi). \quad (36)$$

*Proof.* Without loss of generality we assume that  $j = 1$ .

Set  $H_{2^k-1}(t) := \sigma(t - 2^{-k}s)/(1 - 2^{-k}s)$ ,  $H_{2^k}(t) := 0$  and  $H_s := \varphi_{\bar{\mathbf{k}}, s+1}^*$  for  $s \in Z(k) \setminus \{2^k - 1\}$ , where  $\varphi_{\bar{\mathbf{k}}, s+1}^*$  is defined as in (6). Let

$$Z_i(k) := \{s \in Z(k) : s = 3r + i, r \in \mathbb{N}_0\}, \quad i = 0, 1, 2.$$

To make the proof simple, we divide it into several cases of  $\Lambda$  and  $\Phi(\cdot)$ .

*Case 1.* The case  $\Lambda \subset Z_i(k)$  for some  $i \in \{0, 1, 2\}$  and  $\Phi(\mathbf{x}) \geq 0$  for all  $\mathbf{x} \in \mathbb{I}^d$ . We will show that

$$f(\mathbf{x}) = \sigma \left( \Phi \left( \sum_{s \in \Lambda} H_s(x_1), \bar{\mathbf{x}}_1 \right) - \Phi \left( 1 - \sum_{s \in \Lambda} H_{s+1}(x_1), \bar{\mathbf{x}}_1 \right) \right), \quad (37)$$

for all  $\mathbf{x} \in \mathbb{I}^d$ . Indeed, if  $x_1 \notin \cup_{s \in \Lambda} [2^{-k}s, 2^{-k}(s+3)]$  we have  $\sum_{s \in \Lambda} H_{s+1}(x_1) = \sum_{s \in \Lambda} H_s(x_1) = 0$ . Since  $\text{supp } \Phi(\cdot) \subset \mathbb{I}^d$  we get

$$f(\mathbf{x}) = 0 = \sigma(\Phi(0, \bar{\mathbf{x}}_1) - \Phi(1, \bar{\mathbf{x}}_1)).$$

If  $x_1 \in [2^{-k}s_0, 2^{-k}(s_0 + 1)]$  for some  $s_0 \in \Lambda$  we have  $\sum_{s \in \Lambda} H_{s+1}(x_1) = 0$  and  $\sum_{s \in \Lambda} H_s(x_1) = 2^k x_1 - s_0$ . Since  $\Phi(\mathbf{x}) \geq 0$  and  $\text{supp } \Phi(\cdot) \subset \mathbb{I}^d$  we obtain

$$f(\mathbf{x}) = \Phi(2^k x_1 - s_0, \bar{\mathbf{x}}_1) = \sigma(\Phi(2^k x_1 - s_0, \bar{\mathbf{x}}_1) - \Phi(1, \bar{\mathbf{x}}_1)).$$

If  $x_1 \in [2^{-k}(s_0 + 2), 2^{-k}(s_0 + 3)]$  for some  $s_0 \in \Lambda$  we have  $\sum_{s \in \Lambda} H_s(x_1) = 0$ . Again from  $\Phi(\mathbf{x}) \geq 0$  and  $\text{supp } \Phi(\cdot) \subset \mathbb{I}^d$  we get

$$f(\mathbf{x}) = 0 = \sigma\left(\Phi(0, \bar{\mathbf{x}}_1) - \Phi\left(1 - \sum_{s \in \Lambda} H_{s+1}(x_1), \bar{\mathbf{x}}_1\right)\right).$$

If  $x_1 \in [2^{-k}(s_0 + 1), 2^{-k}(s_0 + 2)]$ ,  $s_0 \in \Lambda$ , it is easy to see that  $\sum_{s \in \Lambda} H_s(x_1) = 1 - \sum_{s \in \Lambda} H_{s+1}(x_1)$ . Hence, the equality (37) holds. We have

$$H_s(x_1) = \sigma(1 - \sigma(2^k x_1 - s - 1) - \sigma(s + 1 - 2^k x_1))$$

for  $s \in Z(k) \setminus \{2^k - 1\}$  and  $H_{2^k-1}(x_1) = \frac{1}{1-2^{-k}s} \sigma(x_1 - 2^{-k}s)$ .

Denote the neural networks on the right side by  $\Phi(H_s)$ . Then the functions  $\sum_{s \in \Lambda} H_s(x_1)$  and  $1 - \sum_{s \in \Lambda} H_{s+1}(x_1)$  can be realized exactly by two networks  $\Phi_1$  and  $\Phi_2$  constructed by parallelization of  $\Phi_{H_s}$ . By Lemma 1, the length of  $\Phi_1$  and  $\Phi_2$  is 3 and their sizes are bounded  $C|\Lambda|$ . Since  $\Phi_1(x_1) \geq 0$  and  $\Phi_2(x_1) \geq 0$  when  $x_1 \in \mathbb{I}$ , we can write

$$f(\mathbf{x}) = \sigma[\Phi(\sigma(\Phi_1(x_1)), \sigma(\sigma(\bar{\mathbf{x}}_1)))] - \Phi(\sigma(\Phi_2(x_1)), \sigma(\sigma(\bar{\mathbf{x}}_1)))].$$

Therefore, the network  $\Phi_\Lambda$  is a concatenation of  $\Phi_1$ ,  $\Phi_2$ ,  $\sigma(\sigma(\bar{\mathbf{x}}_1))$ , and  $\Phi$ . It is clear that we have the estimate

$$W(\Phi_\Lambda) \leq C(d|\Lambda| + W(\Phi)), \quad L(\Phi_\Lambda) \leq 4 + L(\Phi).$$

*Case 2.* The case  $\Lambda \subset Z_i(\mathbf{k})$  for some  $i \in \{0, 1, 2\}$  and  $\Phi(\mathbf{x})$  changing sign when  $\mathbf{x} \in \mathbb{I}^d$ . In this case, we write  $\Phi(\mathbf{x}) = \sigma(\Phi(\mathbf{x})) - \sigma(-\Phi(\mathbf{x}))$ . Hence

$$f(\mathbf{x}) := \sum_{s \in \Lambda} \sigma(\Phi(2^k x_1 - s, \bar{\mathbf{x}}_1)) - \sum_{s \in \Lambda} \sigma(-\Phi(2^k x_1 - s, \bar{\mathbf{x}}_1)), \quad \mathbf{x} \in \mathbb{I}^d.$$

Applying the construction in Case 1 for each sum on the right side with  $\Phi$  replaced by  $\text{Id}(\sigma(\Phi(\cdot)))$  and  $\text{Id}(\sigma(-\Phi(\cdot)))$  respectively we obtain two neural networks  $\Phi_\Lambda^+$  and  $\Phi_\Lambda^-$ . Here  $\text{Id}$  is the identity operator. Concatenating these two network by parallelization, see Lemma 1, we obtain  $\Phi_\Lambda$ . Note that

$$W(\text{Id}(\sigma(\Phi(\cdot)))) = W(\text{Id}(\sigma(-\Phi(\cdot)))) = W(\Phi) + 1$$

and

$$L(\text{Id}(\sigma(\Phi(\cdot)))) = L(\text{Id}(\sigma(-\Phi(\cdot)))) = L(\Phi) + 1.$$

Therefore, the estimates (36) still hold true.

*Case 3.* General case. We rewrite  $f$  in the form:

$$f(\mathbf{x}) = \sum_{j=0,1,2} \sum_{s \in \Lambda \cap Z_j(k)} \Phi(2^k x_1 - s, \bar{\mathbf{x}}_1).$$

To construct the network  $\Phi_\Lambda$ , we first construct the network  $\Phi_{\Lambda_j}$ ,  $j = 0, 1, 2$ , by using the procedure in Case 2 to have that

$$\Phi_{\Lambda_j}(\mathbf{x}) = \sum_{s \in \Lambda \cap Z_j(k)} \Phi(2^k x_1 - s, \bar{\mathbf{x}}_1).$$

Then by parallelizing  $(\Phi_{\Lambda_j})_{j=0,1,2}$  we obtain the network  $\Phi_\Lambda$ . From Lemma 1 we prove (36).  $\blacksquare$

**Lemma 14.** *Let  $d, m \in \mathbb{N}$ ,  $d \geq 2$ ,  $\mathbf{k} \in \mathbb{N}^d$ ,  $\alpha \in (0, 1]$  and  $\varepsilon \in (0, 1)$ . Assume that  $\Phi_\varepsilon(S)$  is the neural network constructed in Lemma 12 to approximate  $S \in \mathcal{S}^{\alpha, d}(m)$  with accuracy  $\varepsilon$  and computation complexity as in (27) and (28), (29). Then for every  $f \in \mathring{U}_\infty^{\alpha, d}$  we can explicitly construct a deep ReLU neural network  $\Phi_\varepsilon(S_{\mathbf{k}, m}(f))$  so that*

$$\|\Phi_\varepsilon(S_{\mathbf{k}, m}(f)) - S_{\mathbf{k}, m}(f)\|_\infty \leq 2^{-\alpha|\mathbf{k}|_1 + d} \varepsilon. \quad (38)$$

Moreover,

$$W(\Phi_\varepsilon(S_{\mathbf{k}, m}(f))) \leq Cd \left( 2^{|\mathbf{k}|_1} + \log d 2^{|\mathbf{k}|_1 - |\mathbf{k}|_\infty} N_d(m) 2^m \binom{m+d-1}{d-1} \log(dB^d \varepsilon^{-1}) \right), \quad (39)$$

and

$$L(\Phi_\varepsilon(S_{\mathbf{k}, m}(f))) \leq C \log d N_d(m) 2^m \log(dB^d \varepsilon^{-1}), \quad (40)$$

where  $N_d(m)$  is given in Lemma 6 and  $B$  is given in Lemma 5.

*Proof.* We can assume without loss of generality that  $k_1 = |\mathbf{k}|_\infty$ . By the definition (25), for  $f \in \mathring{U}_\infty^{\alpha, d}$  we have that

$$S_{\mathbf{k}, m}(f)(\mathbf{x}) := 2^{-\alpha|\mathbf{k}|_1 + d} \sum_{\bar{\mathbf{s}}_1 \in Z(\bar{\mathbf{k}}_1)} \sum_{s_1 \in Z(k_1)} S_m(T_{\mathbf{k}, s}(f))(2^k \mathbf{x} - s).$$

We number the elements of the set  $\mathcal{S}^{\alpha, d}(m)$  from 1 to  $N_d(m)$  as  $S_1, \dots, S_{N_d(m)}$ . For  $\bar{\mathbf{s}}_1 \in Z(\bar{\mathbf{k}}_1)$  and  $\eta = 1, \dots, N_d(m)$ , we define

$$\Lambda_\eta(\bar{\mathbf{s}}_1) := \{s_1 \in Z(k_1) : S_m(T_{\mathbf{k}, s}(f)) = S_\eta \in \mathcal{S}^{\alpha, d}(m)\}.$$

Hence, we can write

$$S_{\mathbf{k}, m}(f)(\mathbf{x}) = 2^{-\alpha|\mathbf{k}|_1 + d} \sum_{\bar{\mathbf{s}}_1 \in Z(\bar{\mathbf{k}}_1)} \sum_{\eta=1}^{N_d(m)} \sum_{s_1 \in \Lambda_\eta(\bar{\mathbf{s}}_1)} S_\eta(2^k \mathbf{x} - s).$$

To approximate  $S_{\mathbf{k}, m}(f)$  we use the output

$$\Phi_\varepsilon(S_{\mathbf{k}, m}(f))(\mathbf{x}) := 2^{-\alpha|\mathbf{k}|_1 + d} \sum_{\bar{\mathbf{s}}_1 \in Z(\bar{\mathbf{k}}_1)} \sum_{\eta=1}^{N_d(m)} \sum_{s_1 \in \Lambda_\eta(\bar{\mathbf{s}}_1)} \Phi_\varepsilon(S_\eta)(2^k \mathbf{x} - s) \quad (41)$$

of a deep ReLU neural network  $\Phi_\varepsilon(S_{\mathbf{k},m}(f))$ . Let us first show explicitly how to construct such a network  $\Phi_\varepsilon(S_{\mathbf{k},m}(f))$  and then estimate its size and depth. Denote by  $\Phi_{S_\eta, \bar{s}_1}$  the network constructed by adding a layer of  $d$  nodes before the input layer of  $\Phi_\varepsilon(S_\eta)$ . Computations at nodes in the first layer of  $\Phi_{S_\eta, \bar{s}_1}$  are  $\sigma(x_1)$  and  $\sigma(2^{k_j}x_j - s_j)$ ,  $j = 2, \dots, d$ . Then by (28) and (29) we have

$$W(\Phi_{S_\eta, \bar{s}_1}) \leq 2(d-1) + 1 + W(\Phi_\varepsilon(S_\eta)) \leq CW(\Phi_\varepsilon(S_\eta))$$

and

$$L(\Phi_{S_\eta, \bar{s}_1}) \leq 1 + L(\Phi_\varepsilon(S_\eta)).$$

Since  $\text{supp } \Phi_\varepsilon(S_\eta) \subset \mathbb{I}^d$ , we have  $\Phi_{S_\eta, \bar{s}_1}(\mathbf{x}) = \Phi_\varepsilon(S_\eta)(x_1, 2^{\bar{k}_1}\bar{x}_1 - \bar{s}_1)$ . Hence we can write

$$\Phi_\varepsilon(S_{\mathbf{k},m}(f))(\mathbf{x}) = 2^{d-\alpha|\mathbf{k}|_1} \sum_{\bar{s}_1 \in Z(\bar{\mathbf{k}}_1)} \sum_{\eta=1}^{N_d(m)} \sum_{s_1 \in \Lambda_\eta(\bar{s}_1)} \Phi_{S_\eta, \bar{s}_1}(2^{k_1}x_1 - s_1, \bar{\mathbf{x}}_1).$$

Applying Lemma 13 to the function  $\sum_{s_1 \in \Lambda_\eta(\bar{s}_1)} \Phi_{S_\eta, \bar{s}_1}(2^{k_1}x_1 - s_1, \bar{\mathbf{x}}_1)$ , we can explicitly construct a network  $\Phi_{\Lambda_\eta(\bar{s}_1)}$  with the output

$$\Phi_{\Lambda_\eta(\bar{s}_1)}(\mathbf{x}) = \sum_{s_1 \in \Lambda_\eta(\bar{s}_1)} \Phi_{S_\eta, \bar{s}_1}(2^{k_1}x_1 - s_1, \bar{\mathbf{x}}_1),$$

so that its size and depth satisfy

$$W(\Phi_{\Lambda_\eta(\bar{s}_1)}) \leq C(d|\Lambda_\eta(\bar{s}_1)| + W(\Phi_{S_\eta, \bar{s}_1})) \leq C(d|\Lambda_\eta(\bar{s}_1)| + W(\Phi_\varepsilon(S_\eta)))$$

and

$$L(\Phi_{\Lambda_\eta(\bar{s}_1)}) \leq CL(\Phi_{S_\eta, \bar{s}_1}) \leq CL(\Phi_\varepsilon(S_\eta)).$$

Let  $\Phi_{\bar{s}_1}$  be the special network combining  $(\Phi_{\Lambda_\eta(\bar{s}_1)})_{\eta=1, \dots, N_d(m)}$  with output

$$\Phi_{\bar{s}_1}(\mathbf{x}) = \sum_{\eta=1}^{N_d(m)} \Phi_{\Lambda_\eta(\bar{s}_1)}(\mathbf{x}).$$

By Lemmata 2 and 12 its length is bounded as

$$L(\Phi_{\bar{s}_1}) \leq \sum_{\eta=1}^{N_d(m)} L(\Phi_{\Lambda_\eta(\bar{s}_1)}) \leq C \sum_{\eta=1}^{N_d(m)} L(\Phi_\varepsilon(S_\eta)) \leq C \log d N_d(m) 2^m \log(dB^d \varepsilon^{-1})$$

and its size

$$\begin{aligned} W(\Phi_{\bar{s}_1}) &\leq \sum_{\eta=1}^{N_d(m)} W(\Phi_{\Lambda_\eta(\bar{s}_1)}) + (d+1)L(\Phi_{\bar{s}_1}) \\ &\leq \sum_{\eta=1}^{N_d(m)} C(d|\Lambda_\eta(\bar{s}_1)| + W(\Phi_\varepsilon(S_\eta))) + (d+1)L(\Phi_{\bar{s}_1}) \\ &\leq C \left( \sum_{\eta=1}^{N_d(m)} d|\Lambda_\eta(\bar{s}_1)| + (d \log d) N_d(m) 2^m \log(dB^d \varepsilon^{-1}) \binom{m+d-1}{d-1} \right) \\ &\leq Cd \left( 2^{k_1} + (\log d) N_d(m) 2^m \binom{m+d-1}{d-1} \log(dB^d \varepsilon^{-1}) \right). \end{aligned}$$

Since  $\mathbf{x} \in \mathbb{I}^d$ , the network  $\Phi_{\bar{\mathbf{s}}_1}$  can be transformed to a standard ReLU neural network with the same output and estimation for depth and size (by adjusting the constants), see (3) and (4). We still denote this new network by  $\Phi_{\bar{\mathbf{s}}_1}$ .

The network  $\Phi_\varepsilon(S_{\mathbf{k},m}(f))$  is a parallelization of  $(\Phi_{\bar{\mathbf{s}}_1})_{\bar{\mathbf{s}}_1 \in \bar{\mathbf{k}}_1}$  which has output (41) and by Lemma 1

$$\begin{aligned} W(\Phi_\varepsilon(S_{\mathbf{k},m}(f))) &\leq 2 \cdot 2^{|\bar{\mathbf{k}}_1|_1} \max_{\bar{\mathbf{s}}_1 \in \bar{\mathbf{k}}_1} W(\Phi_{\bar{\mathbf{s}}_1}) \\ &\leq Cd2^{|\bar{\mathbf{k}}_1|_1} \left( 2^{k_1} + \log dN_d(m)2^m \binom{m+d-1}{d-1} \log(dB^d\varepsilon^{-1}) \right) \end{aligned}$$

and

$$L(\Phi_\varepsilon(S_{\mathbf{k},m}(f))) \leq \max_{\bar{\mathbf{s}}_1 \in Z(\bar{\mathbf{k}}_1)} L(\Phi_{\bar{\mathbf{s}}_1}) \leq C \log dN_d(m)2^m \log(dB^d\varepsilon^{-1}).$$

Thus, (39) and (40) have been proven. Next, we prove the estimate of the approximation error (38). Notice that by the assumptions of the lemma and Lemma 12  $\text{supp } S \subset \mathbb{I}^d$  and  $\text{supp } \Phi_\varepsilon(S) \subset \mathbb{I}^d$  for all  $S \in \mathcal{S}^{\alpha,d}(m)$ , and it holds the estimate (27). Moreover, for different pairs  $(\mathbf{s}, \eta)$  and  $(\mathbf{s}', \eta')$ , the supports of the functions  $(\Phi_\varepsilon(S_\eta) - S_\eta)(2^{\mathbf{k}} \cdot -\mathbf{s})$  and  $(\Phi_\varepsilon(S_{\eta'}) - S_{\eta'})(2^{\mathbf{k}} \cdot -\mathbf{s}')$  are disjoint. Hence, by (27) we obtain

$$\begin{aligned} \|\Phi_\varepsilon(S_{\mathbf{k},m}(f)) - S_{\mathbf{k},m}(f)\|_\infty &= 2^{d-\alpha|\mathbf{k}|_1} \left\| \sum_{\bar{\mathbf{s}}_1 \in Z(\bar{\mathbf{k}}_1)} \sum_{\eta=1}^{N_d(m)} \sum_{\mathbf{s}_1 \in \Lambda_\eta(\bar{\mathbf{s}}_1)} (\Phi_\varepsilon(S_\eta) - S_\eta)(2^{\mathbf{k}} \cdot -\mathbf{s}) \right\|_\infty \\ &= 2^{d-\alpha|\mathbf{k}|_1} \max_{\bar{\mathbf{s}}_1 \in Z(\bar{\mathbf{k}}_1)} \max_{1 \leq \eta \leq N_d(m)} \max_{\mathbf{s}_1 \in \Lambda_\eta(\bar{\mathbf{s}}_1)} \|(\Phi_\varepsilon(S_\eta) - S_\eta)(2^{\mathbf{k}} \cdot -\mathbf{s})\|_\infty \\ &\leq 2^{d-\alpha|\mathbf{k}|_1} \varepsilon \end{aligned}$$

which proves (38). ■

We are now in position to prove Theorem 11.

*Proof.* For convenience, we divide the proof into several steps.

*Step 1.* Let us recall our plan of the proof. To approximate  $f \in \mathring{U}_\infty^{\alpha,d}$ , we will construct a deep ReLU neural network with an output of the form

$$\Phi_\varepsilon(f) = \Phi_{\varepsilon/2}(R_n(f)) + \Phi_{\varepsilon/2}(f - R_n(f)), \quad (42)$$

where  $\Phi_{\varepsilon/2}(R_n(f))$  and  $\Phi_{\varepsilon/2}(f - R_n(f))$  are deep ReLU neural networks approximating  $R_n(f)$  and  $f - R_n(f)$  with accuracy  $\varepsilon/2$ , respectively. Then we have

$$\|f - \Phi_\varepsilon(f)\|_\infty \leq \|R_n(f) - \Phi_{\varepsilon/2}(R_n(f))\|_\infty + \|(f - R_n(f)) - \Phi_{\varepsilon/2}(f - R_n(f))\|_\infty \leq \varepsilon. \quad (43)$$

For approximation of the first term  $R_n(f)$ , we take the deep ReLU neural network  $\Phi_{\varepsilon/2}(R_n(f))$  which has been constructed in Lemma 10. In the following we construct a deep ReLU neural network  $\Phi_{\varepsilon/2}(f - R_n(f))$  for approximating  $f - R_n(f)$  with accuracy  $\varepsilon/2$ .

As noticed above, since the difference  $f - R_n(f)$  is represented as in (23), we shall explicitly construct deep ReLU neural networks  $\Phi_{\varepsilon'}(F_{\mathbf{k}_j})$  to approximate each term  $F_{\mathbf{k}_j}$  with accuracy  $\varepsilon'$  in the sum in (23), where the value of  $\varepsilon'$  will be chosen latter. For ease of

notation we consider the case  $\text{supp}(\mathbf{k}_j) = j$  with  $1 \leq j \leq d-1$ . The other cases are carried out similarly with a slight modification. From (24) we have

$$\begin{aligned} F_{\mathbf{k}_j} &= \prod_{i=1}^j (T_{(k_i-1)\mathbf{e}^j} - T_{k_i\mathbf{e}^j}) T_{(n+1-|\mathbf{k}_j|_1)\mathbf{e}^{j+1}}(f) \\ &= \sum_{\mathbf{e} \in \{0,1\}^j} (-1)^{|\mathbf{e}|_1} T_{\mathbf{k}_j - \mathbf{e}} T_{(n+1-|\mathbf{k}_j|_1)\mathbf{e}^{j+1}}(f) = \sum_{\boldsymbol{\ell} \in \Lambda(\mathbf{k}_j)} c_{\boldsymbol{\ell}} T_{\boldsymbol{\ell}}(f), \end{aligned}$$

where

$$\Lambda(\mathbf{k}_j) := \left\{ \boldsymbol{\ell} \in \mathbb{N}_0^d, \text{supp}(\boldsymbol{\ell}) \subset \{1, \dots, j+1\}, \ell_j = \mathbf{k}_j - \mathbf{e}, \ell_{j+1} = n+1 - |\mathbf{k}_j|_1, \mathbf{e} \in \{0,1\}^j \right\}$$

and  $c_{\boldsymbol{\ell}}$  is either 1 or  $-1$ . It is easy to see that  $|\Lambda(\mathbf{k}_j)| \leq 2^j$  for all  $\mathbf{k}_j$  and if  $\boldsymbol{\ell} \in \Lambda(\mathbf{k}_j)$  then  $n+1-d \leq |\boldsymbol{\ell}|_1 \leq n+1$ .

We approximate  $F_{\mathbf{k}_j}$  by the output

$$\Phi_{\varepsilon'}(F_{\mathbf{k}_j})(\mathbf{x}) := \sum_{\boldsymbol{\ell} \in \Lambda(\mathbf{k}_j)} c_{\boldsymbol{\ell}} \Phi_{\varepsilon'}(S_{\boldsymbol{\ell},m}(f)),$$

where the networks  $\Phi_{\varepsilon'}(S_{\boldsymbol{\ell},m}(f))$  are constructed as in Lemma 14. The network  $\Phi_{\varepsilon'}(F_{\mathbf{k}_j})$  is a parallelization of  $\Phi_{\varepsilon'}(S_{\boldsymbol{\ell},m}(f))$ ,  $\boldsymbol{\ell} \in \Lambda(\mathbf{k}_j)$ .

We define  $\Phi_{\varepsilon/2}(f - R_n(f))$  as a deep ReLU neural network with the output

$$\Phi_{\varepsilon/2}(f - R_n(f)) := \sum_{j=0}^{d-1} \sum_{|\mathbf{k}_j|_1 \leq n} \Phi_{\varepsilon'}(F_{\mathbf{k}_j})(\mathbf{x}) \quad (44)$$

which is a parallelization of  $\Phi_{\varepsilon'}(F_{\mathbf{k}_j})$ ,  $|\mathbf{k}_j|_1 \leq n$ ,  $j = 0, \dots, d-1$ . It approximates  $f - R_n(f)$  with accuracy  $\varepsilon/2$  by an appropriate choice of  $\varepsilon'$ .

We put

$$\varepsilon' = B^d 2^{-\alpha m} \binom{m+d}{d-1}$$

with  $m$  will be chosen later such that  $\varepsilon' \in (0,1)$ . We have from (26) and Lemma 14

$$\begin{aligned} \|F_{\mathbf{k}_j} - \Phi_{\varepsilon'}(F_{\mathbf{k}_j})\|_{\infty} &\leq \sum_{\boldsymbol{\ell} \in \Lambda(\mathbf{k}_j)} \left( \|T_{\boldsymbol{\ell}}(f) - S_{\boldsymbol{\ell},m}(f)\|_{\infty} + \|\Phi_{\varepsilon'}(S_{\boldsymbol{\ell},m}(f)) - S_{\boldsymbol{\ell},m}(f)\|_{\infty} \right) \\ &\leq \sum_{\boldsymbol{\ell} \in \Lambda(\mathbf{k}_j)} \left( (2B)^d (2^m 2^{|\boldsymbol{\ell}|_1})^{-\alpha} \binom{m+d}{d-1} + 2^{-\alpha|\boldsymbol{\ell}|_1+d}\varepsilon' \right) \\ &\leq \sum_{\boldsymbol{\ell} \in \Lambda(\mathbf{k}_j)} \left( (2B)^d (2^m 2^{n+1-d})^{-\alpha} \binom{m+d}{d-1} + 2^{-\alpha(n+1-d)+d}\varepsilon' \right) \\ &\leq 2^{j+1-\alpha} (2^{\alpha+1} B)^d (2^m 2^n)^{-\alpha} \binom{m+d}{d-1}. \end{aligned}$$

This leads to

$$\begin{aligned}
\|(f - R_n(f)) - \Phi_{\varepsilon/2}(f - R_n(f))\|_\infty &\leq \sum_{j=0}^{d-1} \sum_{|\mathbf{k}_j|_1 \leq n} \|\Phi_{\varepsilon'}(F_{\mathbf{k}_j}) - F_{\mathbf{k}_j}\|_\infty \\
&\leq \sum_{j=0}^{d-1} \sum_{|\mathbf{k}_j|_1 \leq n} 2^{j+1-\alpha} (2^{\alpha+1}B)^d (2^m 2^n)^{-\alpha} \binom{m+d}{d-1} \\
&\leq \sum_{j=0}^{d-1} 2^j \binom{n+j}{j} 2^{1-\alpha} (2^{\alpha+1}B)^d (2^m 2^n)^{-\alpha} \binom{m+d}{d-1} \\
&\leq 2^{1-\alpha} (2^{\alpha+2}B)^d (2^m 2^n)^{-\alpha} \binom{m+d}{d-1} \binom{n+d-1}{d-1}.
\end{aligned}$$

We denote the last term by  $A_{n,m}$ . In the next step, our task is to choose  $n, m$  (and therefore,  $\varepsilon'$ ) depending on  $\varepsilon$  such that  $A_{n,m} \leq \varepsilon/2$ . Then we define the deep ReLU neural network  $\Phi_\varepsilon(f)$  as a parallelization of the networks  $\Phi_{\varepsilon/2}(R_n(f))$  and  $\Phi_{\varepsilon/2}(f - R_n(f))$  with the output (42). From this (43) follows. The size and depth of  $\Phi_\varepsilon(f)$  are estimated explicitly in  $d$  and  $\varepsilon$  from the estimation of sizes and depths of  $\Phi_{\varepsilon/2}(R_n(f))$  and  $\Phi_{\varepsilon'}(F_{\mathbf{k}_j})$  by the choice of  $m, n$ .

*Step 2. The choices of  $\varepsilon_0$  and  $n, m$ .* Define  $m_0 \geq d$  as the smallest integer such that  $B^d 2^{-\alpha m_0} \binom{m_0+d}{d-1} < 1$ . Denote  $n_0 \in \mathbb{N}$  from which the function

$$h(n) := K_{d,\alpha} 2^{-\alpha n} n^{d-1-\alpha} (\log n)^{(\alpha+1)(d-1)}, \quad (45)$$

where

$$K_{d,\alpha} := 2(2^{\alpha+2}B)^d (4d \log 3)^\alpha \left( \frac{2^{d-1}}{(d-1)!} \right)^{\alpha+2} \quad (46)$$

is decreasing and  $h(n-1) \leq 2^{-\alpha n/2}$  for all  $n \geq n_0$ . We put  $n_1 = \lfloor (8d \log 3) 2^{m_0} \binom{m_0+d-1}{d-1} \rfloor + 1$  and define  $\varepsilon_0 = \min\{h(n_0), h(n_1), 1/2\}$ . For  $\varepsilon \in (0, \varepsilon_0)$  we choose  $n \in \mathbb{N}$ ,  $n \geq \max\{n_0, n_1\}$ , such that  $h(n) \leq \varepsilon/2 < h(n-1)$  and then  $m$  such that

$$(\log d) 3^{2^{m+1} \binom{m+d-1}{d-1}} 2^m \binom{m+d-1}{d-1} m \leq 2^{\frac{n}{d}} < (\log d) 3^{2^{m+2} \binom{m+d}{d-1}} 2^{m+1} \binom{m+d}{d-1} (m+1). \quad (47)$$

These choices imply

$$3^{2^{m+1} \binom{m+d-1}{d-1}} \leq 2^{\frac{n}{d}} < 3^{2^{m+3} \binom{m+d}{d-1}}$$

and

$$2^{m+1} \binom{m+d-1}{d-1} \log 3 < \frac{n}{d} < (8 \log 3) 2^m \binom{m+d}{d-1} \quad \text{and} \quad m \leq \log n. \quad (48)$$

Since  $n \geq (8d \log 3) 2^{m_0} \binom{m_0+d}{d-1}$  we get  $m \geq m_0 \geq d$  and

$$\begin{aligned}
A_{n,m} &\leq 2^{1-\alpha} (2^{\alpha+2}B)^d 2^{-n\alpha} \left[ dn^{-1} (8 \log 3) \binom{m+d}{d-1} \right]^\alpha \binom{m+d}{d-1} \binom{n+d-1}{d-1} \\
&\leq 2(2^{\alpha+2}B)^d (4d \log 3)^\alpha \left( \frac{2^{d-1}}{(d-1)!} \right)^{\alpha+2} 2^{-\alpha n} n^{d-1-\alpha} m^{(\alpha+1)(d-1)} \\
&= h(n) \leq \frac{\varepsilon}{2}.
\end{aligned}$$

*Step 3. Estimating the size and depth of  $\Phi_{\varepsilon'}(S_{\ell,m}(f))$ .* From  $n+1-d \leq |\ell|_1 \leq n+1$  we have



$$|\ell|_1 - |\ell|_\infty \leq n + 1 - \frac{n + 1 - d}{d} \leq n - \frac{n}{d} + 2$$

which by (39) leads to

$$W(\Phi_{\varepsilon'}(S_{\ell,m}(f))) \leq Cd \left( 2^n + (\log d) 2^{n - \frac{n}{d} + m} \binom{m + d - 1}{d - 1} \log(dB^d \varepsilon'^{-1}) N_d(m) \right).$$

Note that by the choice of  $\varepsilon'$  we get

$$\log(dB^d \varepsilon'^{-1}) \leq \log \left( d 2^{\alpha m} \binom{m + d}{d - 1}^{-1} \right) \leq \alpha m.$$

It yields from (47)

$$(\log d) 2^m \binom{m + d - 1}{d - 1} m N_d(m) \leq 3^{2^{m+1} \binom{m+d-1}{d-1}} 2^m \binom{m + d - 1}{d - 1} m \leq 2^{\frac{n}{d}}.$$

Consequently

$$W(\Phi_{\varepsilon'}(S_{\ell,m}(f))) \leq C_\alpha d 2^n. \quad (49)$$

Similarly, we have

$$L(\Phi_{\varepsilon'}(S_{\ell,m}(f))) \leq C \log d N_d(m) 2^m \log(dB^d \varepsilon'^{-1}) \leq C_\alpha (\log d) 3^{2^{m+1} \binom{m+d-1}{d-1}} 2^m m \leq C_\alpha 2^{\frac{n}{d}}.$$

*Step 4. Estimation of the size and depth of  $\Phi_\varepsilon(f)$ .* We recall that  $\Phi_{\varepsilon/2}(f - R_n(f))$  is the network obtained by parallelization of  $\Phi_{\varepsilon'}(S_{\ell,m}(f))$  with  $\ell$  in the multi-set

$$\Lambda = \{ \ell \in \Lambda(\mathbf{k}_j), j = 0, \dots, d - 1, |\mathbf{k}_j|_1 \leq n \}$$

and has the output equal to the double sum on the right side of (44). We have

$$|\Lambda| \leq \sum_{j=0}^{d-1} \sum_{|\mathbf{k}_j|_1 \leq n} 2^j = \sum_{j=0}^{d-1} 2^j \binom{n + j}{j} \leq 2^d \binom{n + d - 1}{d - 1}.$$

The network  $\Phi_\varepsilon(f)$  is a parallelization of  $\Phi_{\varepsilon/2}(R_n(f))$  and  $\Phi_{\varepsilon/2}(f - R_n(f))$ . Therefore, by Lemma 1 and the construction of  $\Phi_{\varepsilon/2}(R_n(f))$  and  $\Phi_{\varepsilon/2}(f - R_n(f))$  we obtain

$$\begin{aligned} W(\Phi_\varepsilon(f)) &\leq C \max \{ W(\Phi_{\varepsilon/2}(R_n(f))), W(\Phi_{\varepsilon/2}(f - R_n(f))) \} \\ &\leq C \max \left\{ W(\Phi_{\varepsilon/2}(R_n(f))), 2^d \binom{n + d - 1}{d - 1} \max_{\ell \in \Lambda} W(\Phi_{\varepsilon'}(S_{\ell,m}(f))) \right\}. \end{aligned}$$

From Lemma 10 and (49) we deduce that

$$W(\Phi_\varepsilon(f)) \leq C_\alpha \max \left\{ d 2^n \log(dB^d 2\varepsilon^{-1}) \binom{n + d - 1}{d - 1}, d 2^d 2^n \binom{n + d - 1}{d - 1} \right\}.$$

Since  $h(n) \leq \varepsilon/2$ , from (45) by simple calculation we get

$$\log(dB^d \varepsilon^{-1}) \leq C_\alpha (d + n) \leq C_\alpha dn$$

which implies

$$W(\Phi_\varepsilon(f)) \leq C_\alpha d 2^d 2^n n^d \frac{2^{d-1}}{(d-1)!}.$$

From  $\varepsilon/2 \leq h(n - 1) \leq 2^{-\alpha n/2}$  we deduce  $n \leq 2\alpha^{-1} \log(2\varepsilon^{-1})$  and by (45)

$$2^n \leq C_\alpha \left( K_{d,\alpha} 2\varepsilon^{-1} n^{d-1-\alpha} (\log n)^{(\alpha+1)(d-1)} \right)^{\frac{1}{\alpha}}. \quad (50)$$

Consequently

$$\begin{aligned} W(\Phi_\varepsilon(f)) &\leq C_\alpha d 2^d \frac{2^{d-1}}{(d-1)!} \left( K_{d,\alpha} \varepsilon^{-1} (n \log n)^{(\alpha+1)(d-1)} \right)^{\frac{1}{\alpha}} \\ &\leq C_\alpha d 2^d \frac{2^{d-1}}{(d-1)!} (K_{d,\alpha})^{\frac{1}{\alpha}} \varepsilon^{-\frac{1}{\alpha}} \left( \log(2\alpha^{-1} \log(2\varepsilon^{-1})) 2\alpha^{-1} \log(2\varepsilon^{-1}) \right)^{(1+\frac{1}{\alpha})(d-1)}. \end{aligned}$$

We use the inequalities with  $p := (1 + \frac{1}{\alpha})(d-1) \geq 1$ ,  $\varepsilon \in (0, 1/2)$

$$\begin{aligned} [\log(2\alpha^{-1} \log(2\varepsilon^{-1}))]^p &= [\log(2\alpha^{-1}) + \log \log(2\varepsilon^{-1})]^p \\ &\leq [2 \log(2\alpha^{-1}) (\log \log(2\varepsilon^{-1}))]^p \end{aligned} \quad (51)$$

to obtain

$$W(\Phi_\varepsilon(f)) \leq C_\alpha d 2^d \frac{2^{d-1}}{(d-1)!} (K_{d,\alpha})^{\frac{1}{\alpha}} \varepsilon^{-\frac{1}{\alpha}} \left( 4\alpha^{-1} \log(2\alpha^{-1}) \log \log(2\varepsilon^{-1}) \log(2\varepsilon^{-1}) \right)^{(1+\frac{1}{\alpha})(d-1)}.$$

Replacing  $K_{d,\alpha}$  by the right-hand side of (46), we prove (21).

Now we estimate the depth of  $\Phi_\varepsilon(f)$ . By applying Lemmata 2, 10 and by the construction of  $\Phi_{\varepsilon/2}(f - R_n(f))$  we have that

$$\begin{aligned} L(\Phi_\varepsilon(f)) &= C \max \{ L(\Phi_{\varepsilon/2}(R_n(f))), L(\Phi_{\varepsilon/2}(f - R_n(f))) \} \\ &\leq C \max \{ \log d \log(dB^d(\varepsilon/2)^{-1}), \max_{\ell \in \Lambda} L(\Phi_{\varepsilon'}(S_{\ell,m}(f))) \} \\ &\leq C_\alpha \max \{ \log d \log(dB^d \varepsilon^{-1}), 2^{\frac{n}{d}} \} \\ &\leq C_\alpha \max \{ (d \log d)n, 2^{\frac{n}{d}} \}. \end{aligned}$$

By (47) and (48) it is easily seen that  $2^{\frac{n}{d}}$  dominates  $(d \log d)n$ . From (50), the inequality  $n \leq 2\alpha^{-1} \log(2\varepsilon^{-1})$ , and (51) we get

$$\begin{aligned} 2^{\frac{n}{d}} &\leq C_\alpha \left( K_{d,\alpha} \varepsilon^{-1} n^{d-1-\alpha} (\log n)^{(\alpha+1)(d-1)} \right)^{1/(d\alpha)} \\ &\leq C_\alpha \left( K_{d,\alpha} \varepsilon^{-1} (2\alpha^{-1} \log(2\varepsilon^{-1}))^{d-1-\alpha} (\log(2\alpha^{-1} \log(2\varepsilon^{-1})))^{(\alpha+1)(d-1)} \right)^{1/(d\alpha)} \\ &\leq C_\alpha \left( K_{d,\alpha} \varepsilon^{-1} (2\alpha^{-1} \log(2\varepsilon^{-1}))^{d-1-\alpha} (2 \log(2\alpha^{-1}) \log \log(2\varepsilon^{-1}))^{(\alpha+1)(d-1)} \right)^{1/(d\alpha)}. \end{aligned}$$

In view of (46), we find that

$$\left( K_{d,\alpha} (2\alpha^{-1})^{d-1-\alpha} (2 \log(2\alpha^{-1}))^{(\alpha+1)(d-1)} \right)^{1/(d\alpha)}$$

is bounded by a constant depending only on  $\alpha$ . Consequently

$$L(\Phi_\varepsilon(f)) \leq C_\alpha \varepsilon^{-\frac{1}{d\alpha}} (\log(2\varepsilon^{-1}))^{\frac{d-1-\alpha}{d\alpha}} (\log \log(2\varepsilon^{-1}))^{\frac{(\alpha+1)(d-1)}{d\alpha}}$$

which proves (22). ■

## 6. AN APPLICATION TO NUMERICAL SOLVING PDES

In this section, we apply the results on approximation by deep ReLU neural networks in Sections 4 and 5 to numerical approximation to the solution of elliptic PDEs. Before going

into detail, let us analyze the differences between approximation by deep ReLU neural networks and finite element methods. It is well-known that every deep ReLU neural network in  $\mathbb{R}^d$  represents a continuous piecewise linear function defined on a number of polyhedral subdomains and conversely every continuous piecewise linear function in  $\mathbb{R}^d$  can be represented by a deep ReLU. This shows that there exists some deep ReLU neural network which is at least as good as adaptive finite element for solutions of PDEs. However, the approximation to solutions of PDEs in high dimensions by finite element methods is burdened the curse of dimensionality. For some PDEs, it has been shown that deep neural networks are capable of representing solutions without incurring the curse of dimensionality, see, for instance, [15, 20, 21]. For further discussion about relationship between deep ReLU neural networks and finite element methods in numerical solving PDEs, we refer the reader to [17, 28, 30].

Consider a modeled diffusion elliptic equation

$$-\operatorname{div}(a(\mathbf{x})\nabla u(\mathbf{x})) = f(\mathbf{x}) \quad \text{in } \mathbb{I}^d, \quad u|_{\partial\mathbb{I}^d} = 0,$$

with a function  $f$  and a diffusion coefficient  $a$  having sufficient regularity. Denote by  $V := \mathring{W}_2^1(\mathbb{I}^d)$  the energy space. If  $a$  satisfies the ellipticity assumption

$$0 < a_{\min} \leq a(\mathbf{x}) \leq a_{\max} < \infty, \quad \forall \mathbf{x} \in \mathbb{I}^d,$$

by the well-known Lax-Milgram lemma, there exists a unique solution  $u \in V$  in weak form which satisfies the variational equation

$$\int_{\mathbb{I}^d} a(\mathbf{x})\nabla u(\mathbf{x}) \cdot \nabla v(\mathbf{x}) \, d\mathbf{x} = \int_{\mathbb{I}^d} f(\mathbf{x})v(\mathbf{x}) \, d\mathbf{x}, \quad \forall v \in V.$$

We want to approximate the solution  $u$  by deep ReLU neural networks. The approximation error is measured in the norm of  $L_\infty(\mathbb{I}^d)$ . Assume for the modeled case that  $a$  and  $f$  have Hölder-Nikol'skii mixed smoothness 1, i.e.,  $a, f \in H_\infty^1(\mathbb{I}^d)$ . Then, the solution  $u$  has at least mixed derivatives  $\partial^\alpha u$  with  $\alpha \in \mathbb{N}_0^d$ ,  $\max_{j=1,\dots,d} \alpha_j \leq 1$ , belonging to  $L_2(\mathbb{I}^d)$  [12], and therefore, by embedding for function spaces of mixed smoothness, see [32, Theorem 2.4.1],  $u$  belongs to  $\mathring{H}_\infty^{1/2}(\mathbb{I}^d)$ . For simplicity we assume that  $u \in \mathring{U}_\infty^{1/2}$ .

For the nonadaptive approximation, according to Theorem 7, for any  $\varepsilon > 0$  sufficient small one can explicitly construct a deep neural network architecture  $\mathbb{A}_\varepsilon$  independent of  $f$  and  $a$ , and a deep ReLU neural network  $\Phi_\varepsilon(u)$  having the architecture  $\mathbb{A}_\varepsilon$  such that

$$\|u - \Phi_\varepsilon(u)\|_\infty \leq \varepsilon,$$

$$W(\mathbb{A}_\varepsilon) \leq Cd \left( \frac{K_1^d}{(d-1)!} \right)^3 \varepsilon^{-2} \log(2\varepsilon^{-1})^{3(d-1)+1},$$

and

$$L(\mathbb{A}_\varepsilon) \leq C \log d \log(2\varepsilon^{-1}),$$

where  $K_1 := 8(\sqrt{2} + 1)^{3/2}$ .

For the adaptive approximation, according to Theorem 11, for any  $\varepsilon > 0$  sufficient small one can explicitly construct an adaptive deep ReLU neural network  $\Phi_\varepsilon(u)$  so that

$$\|u - \Phi_\varepsilon(u)\|_\infty \leq \varepsilon,$$

$$W(\Phi_\varepsilon(u)) \leq Cd^2 \left( \frac{K_2^d}{(d-1)!} \right)^6 \varepsilon^{-2} (\log(2\varepsilon^{-1}) \log \log(2\varepsilon^{-1}))^{3(d-1)},$$

and

$$L(\Phi_\varepsilon(u)) \leq C' \varepsilon^{-\frac{2}{d}} (\log(2\varepsilon^{-1}))^{\frac{2d-3}{d}} (\log \log(2\varepsilon^{-1}))^{\frac{3(d-1)}{d}},$$

where  $K_2 := 16((2 + \sqrt{2}))^{1/3}$ .

## 7. CONCLUSIONS

We have presented both nonadaptive and adaptive methods for explicit construction of deep ReLU neural network  $\Phi_\varepsilon(f)$  having an output that approximates functions  $f$  in the Hölder-Nikol'skii spaces with an arbitrary prescribed accuracy  $\varepsilon$  in the  $L_\infty$ -norm. Nonadaptivity means that the architecture of approximating deep ReLU neural networks is the same for all functions in  $\mathring{U}_\infty^{\alpha,d}$ . For nonadaptive approximation, by using truncation of Faber series as a intermediate approximation, we have established a dimension-dependent estimate for the computation complexity characterized by the size  $W(\Phi_\varepsilon(f))$  estimated by

$$W(\Phi_\varepsilon(f)) \leq C_\alpha d \left( \frac{K_1^d}{(d-1)!} \right)^{\frac{1}{\alpha}+1} \varepsilon^{-\frac{1}{\alpha}} \log(2\varepsilon^{-1})^{(d-1)(\frac{1}{\alpha}+1)+1},$$

where  $K_1 = B^{1/(\alpha+1)}4\alpha^{-1}$  with  $B = (2^\alpha - 1)^{-1}$ .

Concerning adaptive method, for any  $f \in \mathring{U}_\infty^{\alpha,d}$ , we explicitly construct a deep ReLU neural network  $\Phi_\varepsilon(f)$  of adaptive architecture having the output that approximates  $f$  in the  $L_\infty(\mathbb{I}^d)$ -norm with a prescribed accuracy  $\varepsilon$  and having the size estimated by

$$W(\Phi_\varepsilon(f)) \leq C_\alpha d^2 \left( \frac{K_2^d}{(d-1)!} \right)^{\frac{2}{\alpha}+2} \varepsilon^{-\frac{1}{\alpha}} (\log(2\varepsilon^{-1}) \log \log(2\varepsilon^{-1}))^{(1+\frac{1}{\alpha})(d-1)},$$

where  $K_2 = 4(2^{\alpha+3}B)^{\frac{1}{2\alpha+2}}(\alpha^{-1} \log(2\alpha^{-1}))^{1/2}$ .

Construction of deep ReLU neural networks in the adaptive method is more involved but improves  $\log(2\varepsilon^{-1})$  in the computation complexity of the approximating deep ReLU neural networks compared to the nonadaptive one.

Our theory is illustrated by an application to numerical approximation to the solution of elliptic PDEs.

## ACKNOWLEDGMENT

This research is funded by Vietnam National Foundation for Science and Technology Development (NAFOSTED) under Grant number 102.01-2020.03. A part of this work was done when Dinh Dũng and Van Kien Nguyen were working at the Vietnam Institute for Advanced Study in Mathematics (VIASM). They would like to thank the VIASM for providing a fruitful research environment and working condition.

## REFERENCES

- [1] M. Ali and A. Nouy, "Approximation of smoothness classes by deep ReLU networks," *arXiv:2007.15645*, 2020.
- [2] D. Dũng, "B-spline quasi-interpolant representations and sampling recovery of functions with mixed smoothness," *J. Complexity*, vol. 27, pp. 541–567, 2011.
- [3] D. Dũng and V. K. Nguyen, "Deep ReLU neural networks in high-dimensional approximation," *Neural Netw.*, vol. 142, pp. 619–635, 2021.
- [4] —, "High-dimensional nonlinear approximation by parametric manifolds in Hölder-Nikol'skii spaces of mixed smoothness," *arXiv:2102.04370*, 2021.

- [5] D. Dũng, V. N. Temlyakov, and T. Ullrich, *Hyperbolic Cross Approximation*. Advanced Courses in Mathematics - CRM Barcelona, Birkhäuser/Springer, 2018.
- [6] D. Dũng and M. X. Thao, “Dimension-dependent error estimates for sampling recovery on Smolyak grids based on B-spline quasi-interpolation,” *J. Approx. Theory*, vol. 250, pp. 185–205, 2020.
- [7] I. Daubechies, R. DeVore, S. Foucart, B. Hanin, and G. Petrova, “Nonlinear approximation and (Deep) ReLU networks,” *arXiv:1905.02199*, 2019.
- [8] W. E and Q. Wang, “Exponential convergence of the deep neural network approximation for analytic functions,” *Sci. China Math.*, vol. 61, pp. 1733–1740, 2018.
- [9] M. Geist, P. Petersen, M. Raslan, R. Schneider, and G. Kutyniok, “Numerical solution of the parametric diffusion equation by deep neural networks,” *arXiv:2004.12131*, 2020.
- [10] X. Glorot, A. Bordes, and Y. Bengio, “Deep sparse rectifier neural networks,” *In Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics, Ft. Lauderdale, FL, USA*, pp. 315–323, 2011.
- [11] R. Gribonval, Kutyniok, M. Nielsen, and F. Voigtlaender, “Approximation spaces of deep neural networks,” *arXiv:1905.01208*, 2019.
- [12] M. Griebel and S. Knapek, “Optimized general sparse grid approximation spaces for operator equations,” *Math. Comp.*, vol. 78, pp. 2223–2257, 2009.
- [13] A. Griewank, F. Y. Kuo, H. Leövey, and I. H. Sloan, “High dimensional integration of kinks and jumps – smoothing by preintegration,” *J. Comput. Appl. Math.*, vol. 344, pp. 259–274, 2018.
- [14] P. Grohs, D. Perekrestenko, D. Elbrachter, and H. Bolcskei, “Deep neural network approximation theory,” *arXiv: 1901.02220*, 2019.
- [15] P. Grohs, A. Jentzen, and D. Salimova, “Deep neural network approximations for Monte Carlo algorithms,” *arXiv:1908.10828*, 2019.
- [16] I. Gühring, G. Kutyniok, and P. Petersen, “Error bounds for approximations with deep ReLU neural networks in  $W^{s,p}$  norms,” *Anal. Appl. (Singap.)*, vol. 18, pp. 803–859, 2020.
- [17] J. He, L. Li, J. Xu, and C. Zheng, “Relu deep neural networks and linear finite elements,” *J. Comput. Math*, vol. 38, pp. 502–529, 2020.
- [18] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification,” *2015 IEEE International Conference on Computer Vision*, pp. 1026–1034, 2015.
- [19] D. Hebb, *The Organization of Behavior: A Neuropsychological Theory*. Wiley, 1949.
- [20] M. Hutzenthaler, A. Jentzen, T. Kruse, and T. A. Nguyen, “A proof that rectified deep neural networks overcome the curse of dimensionality in the numerical approximation of semilinear heat equations,” *arXiv:1901.10854*, 2019.
- [21] A. Jentzen, D. Salimova, and T. Welti, “A proof that deep artificial neural networks overcome the curse of dimensionality in the numerical approximation of Kolmogorov partial differential equations with constant diffusion and nonlinear drift coefficients,” *arXiv:1809.07321*, 2018.
- [22] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “ImageNet classification with deep convolutional neural networks,” *NeurIPS*, pp. 1106–1114, 2012.

- [23] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, pp. 436–444, 2015.
- [24] A. L. Maas, A. Y. Hannun, and A. Y. Ng, “Rectifier nonlinearities improve neural network acoustic models,” *Proceedings of the 30 th International Conference on Machine Learning, Atlanta, Georgia, USA*, pp. 315–323, 2014.
- [25] W. S. McCulloch and W. Pitts, “A logical calculus of the ideas immanent in nervous activity,” *Bull. Math. Biophys.*, vol. 5, pp. 115–133, 1943.
- [26] H. N. Mhaskar, “Neural networks for optimal approximation of smooth and analytic functions,” *Neural Comput.*, vol. 8, pp. 164–177, 1996.
- [27] H. Montanelli and Q. Du, “New error bounds for deep ReLU networks using sparse grids,” *SIAM J. Math. Data Sci.*, vol. 1, pp. 78–92, 2019.
- [28] J. A. A. Opschoor, P. C. Petersen, and C. Schwab, “Deep ReLU networks and high-order finite element methods,” *Anal. Appl. (Singap.)*, vol. 18, pp. 715–770, 2020.
- [29] P. Petersen and F. Voigtlaender, “Optimal approximation of piecewise smooth functions using deep ReLU neural networks,” *Neural Netw.*, vol. 108, pp. 296–330, 2018.
- [30] P. C. Petersen, “Neural network theory,” *Preprint*, 2020.
- [31] F. Rosenblatt, “The perceptron: a probabilistic model for information storage and organization in the brain,” *Psychol. Rev.*, vol. 65, pp. 386–408, 1958.
- [32] H. Schmeisser and H. Triebel, *Topics in Fourier Analysis and Function Spaces*. Chichester; New York : Wiley, 1987.
- [33] C. Schwab and J. Zech, “Deep learning in high dimension: Neural network expression rates for generalized polynomial chaos expansions in UQ,” *Anal. Appl. (Singap.)*, vol. 17, pp. 19–55, 2019.
- [34] T. Suzuki, “Adaptivity of deep ReLU network for learning in Besov and mixed smooth Besov spaces: optimal rate and curse of dimensionality,” *International Conference on Learning Representations*, 2019.
- [35] L. Tóth, “Phone recognition with deep sparse rectifier neural networks,” *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 6985–6989, 2013.
- [36] H. Triebel, *Bases in Function Spaces, Sampling, Discrepancy, Numerical Integration*. European Math. Soc. Publishing House, Zürich, 2010.
- [37] —, *Hybrid Function Spaces, Heat and Navier-Stokes Equations*. European Mathematical Society, 2015.
- [38] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, and M. Norouzi, “Google’s neural machine translation system: Bridging the gap between human and machine translation,” *arXiv: 1609.08144*, 2016.
- [39] D. Yarotsky, “Error bounds for approximations with deep ReLU networks,” *Neural Netw.*, vol. 94, pp. 103–114, 2017.
- [40] —, “Quantified advantage of discontinuous weight selection in approximations with deep neural networks,” *arXiv: 1705.01365*, 2017.
- [41] H. Yserentant, *Regularity and Approximability of Electronic Wave Functions*. Lecture Notes in Mathematics, Springer, 2010.

*Received on March 2, 2021*  
*Accepted on August 15, 2021*