

# Neural Network Approximation: Three Hidden Layers Are Enough

Zuowei Shen<sup>a</sup>, Haizhao Yang<sup>b</sup>, Shijun Zhang<sup>a</sup>

<sup>a</sup>*Department of Mathematics, National University of Singapore*

<sup>b</sup>*Department of Mathematics, Purdue University*

---

## Abstract

A three-hidden-layer neural network with super approximation power is introduced. This network is built with the Floor function ( $\lfloor x \rfloor$ ), the exponential function ( $2^x$ ), the step function ( $\mathbb{1}_{x \geq 0}$ ), or their compositions as activation functions in each neuron and hence we call such networks as Floor-Exponential-Step (FLES) networks. For any width hyper-parameter  $N \in \mathbb{N}^+$ , it is shown that FLES networks with a width  $\max\{d, N\}$  and three hidden layers can uniformly approximate a Hölder function  $f$  on  $[0, 1]^d$  with an exponential approximation rate  $3\lambda d^{\alpha/2} 2^{-\alpha N}$ , where  $\alpha \in (0, 1]$  and  $\lambda$  are the Hölder order and constant, respectively. More generally for an arbitrary continuous function  $f$  on  $[0, 1]^d$  with a modulus of continuity  $\omega_f(\cdot)$ , the constructive approximation rate is  $\omega_f(\sqrt{d} 2^{-N}) + 2\omega_f(\sqrt{d}) 2^{-N}$ . As a consequence, this new class of networks overcomes the curse of dimensionality in approximation power when the variation of  $\omega_f(r)$  as  $r \rightarrow 0$  is moderate (e.g.,  $\omega_f(r) \lesssim r^\alpha$  for Hölder continuous functions), since the major term to be concerned in our approximation rate is essentially  $\sqrt{d}$  times a function of  $N$  independent of  $d$  within the modulus of continuity.

*Keywords:*

Exponential Convergence, Curse of Dimensionality, Deep Neural Network, Floor-Exponential-Step Activation Function, Continuous Function.

---

*Email addresses:* [matzuows@nus.edu.sg](mailto:matzuows@nus.edu.sg) (Zuowei Shen), [haizhao@purdue.edu](mailto:haizhao@purdue.edu) (Haizhao Yang), [zhangshijun@u.nus.edu](mailto:zhangshijun@u.nus.edu) (Shijun Zhang)

## 1. Introduction

This paper studies the approximation power of neural networks and shows that three hidden layers are enough for neural networks to achieve super approximation capacity. In particular, leveraging the power of advanced yet simple activation functions, we will introduce new theories and network architectures with only three hidden layers achieving exponential convergence and avoiding the curse of dimensionality simultaneously for (Hölder) continuous functions with an explicit approximation bound. The theories established in this paper would provide new insights to explain why deeper neural networks are better than one-hidden-layer neural networks for large-scale and high-dimensional problems. The approximation theories here are constructive (i.e., with explicit formulas to specify network parameters) and quantitative (i.e., results valid for essentially arbitrary width and/or depth without lower bound constraints) with explicit error bounds working for three-hidden-layer networks with arbitrary width.

Constructive approximation with quantitative results and explicit error bounds would provide important guides for deciding the network sizes in deep learning. For example, the (nearly) optimal approximation rates of deep ReLU networks for a Lipschitz continuous function and a  $C^s$  function  $f$  on  $[0, 1]^d$  are  $\mathcal{O}(\sqrt{d}N^{-2/d}L^{-2/d})$  and  $\mathcal{O}(\|f\|_{C^s}N^{-2s/d}L^{-2s/d})$  Shen et al. (2020); Lu et al. (2020), respectively. For results in terms of the number of nonzero parameters, the reader is referred to Yarotsky (2017); Schmidt-Hieber (2017); Petersen and Voigtlaender (2018); Yarotsky (2018); Gühring et al. (2019); Yarotsky and Zhevnerchuk (2019) and the reference therein. Obviously, the curse of dimensionality exists in ReLU networks for these generic functions and, therefore, ReLU networks would need to be exponentially large in  $d$  to maintain a reasonably good approximation accuracy. The curse could be lessened when target function spaces are smaller. To name a few, Poggio et al. (2017); Barron and Klusowski (2018); E et al. (2019); Montanelli et al. (2020); Chen et al. (2019a); Hutzenthaler et al. (2020) and reference therein for ReLU networks. The limitation of ReLU networks motivated the work in Shen et al. (2020) to introduce Floor-ReLU networks built with either a Floor ( $\lfloor x \rfloor$ ) or ReLU ( $\max\{0, x\}$ ) activation function in each neuron. It was shown by construction in Shen et al. (2020) that Floor-ReLU networks with width  $\max\{d, 5N + 13\}$  and depth  $64dL + 3$  can uniformly approximate a Hölder continuous function  $f$  on  $[0, 1]^d$  with a root-exponential approximation rate  $3\lambda d^{\alpha/2}N^{-\alpha\sqrt{L}}$  without the curse of dimensionality.

38 The most important message of Shen et al. (2020) (and probably also  
 39 of Yarotsky and Zhevnerchuk (2019)) is that the combination of simple acti-  
 40 vation functions can create super approximation power. In the Floor-ReLU  
 41 networks mentioned above, the power of depth is fully reflected in the ap-  
 42 proximation rate  $3\lambda d^{\alpha/2} N^{-\alpha\sqrt{L}}$  that is root-exponential in depth. However,  
 43 the power of width is much weaker and the approximation rate is polyno-  
 44 mial in width if depth is fixed. This seems to be inconsistent with recent  
 45 development of network optimization theory Jacot et al. (2018); Du et al.  
 46 (2019); Mei et al. (2018); Wu et al. (2018); Chen et al. (2019b); Lu et al.  
 47 (2020); Luo and Yang (2020), where larger width instead of depth can ease  
 48 the challenge of highly nonconvex optimization. The mystery of the power  
 49 of width and depth remains and it motivates us to demonstrate that width  
 50 can also enable super approximation power when armed with appropriate  
 51 activation functions. In particular, we explore the Floor function, the expo-  
 52 nential function ( $2^x$ ), the step function ( $\mathbb{1}_{x \geq 0}$ ), or their compositions as activa-  
 53 tion functions to build fully-connected feed-forward neural networks (FNNs).  
 54 These networks are called Floor-Exponential-Step (FLES) networks. As we  
 55 shall prove by construction, Theorem 1.1 below shows that FLES networks  
 56 with width  $\max\{N, d\}$  and three hidden layers can uniformly approximate  
 57 a continuous function  $f$  on  $[0, 1]^d$  with an exponential approximation rate  
 58  $2\omega_f(\sqrt{d})2^{-N} + \omega_f(\sqrt{d}2^{-N})$ , where  $\omega_f(\cdot)$  is the modulus of continuity defined  
 59 as

$$60 \quad \omega_f(r) := \sup \{ |f(\mathbf{x}) - f(\mathbf{y})| : \|\mathbf{x} - \mathbf{y}\|_2 \leq r, \mathbf{x}, \mathbf{y} \in [0, 1]^d \}, \quad \text{for any } r \geq 0,$$

61 where  $\|\mathbf{x}\|_2 = \sqrt{x_1^2 + x_2^2 + \dots + x_d^2}$  for any  $\mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$ . In particular,  
 62 there are three kinds of activation functions denoted as  $\sigma_1$ ,  $\sigma_2$ , and  $\sigma_3$  in  
 63 FLES networks (see Figure 1 for an illustration):

$$64 \quad \sigma_1(x) := \lfloor x \rfloor, \quad \sigma_2(x) := 2^x, \quad \text{and} \quad \sigma_3 := \mathcal{T}(x - \lfloor x \rfloor - \tfrac{1}{2}), \quad \text{for any } x \in \mathbb{R},$$

65 where

$$66 \quad \mathcal{T}(x) := \mathbb{1}_{x \geq 0} = \begin{cases} 1, & x \geq 0, \\ 0, & x < 0, \end{cases} \quad \text{for any } x \in \mathbb{R}.$$

67 **Theorem 1.1.** *Given a continuous function  $f$  on  $[0, 1]^d$  and  $N \in \mathbb{N}^+$ , there*  
 68 *exist  $a_1, a_2, \dots, a_N \in [0, \frac{1}{2})$  such that*

$$69 \quad |\phi(\mathbf{x}) - f(\mathbf{x})| \leq 2\omega_f(\sqrt{d})2^{-N} + \omega_f(\sqrt{d}2^{-N}),$$

70 for any  $\mathbf{x} = (x_1, \dots, x_d) \in [0, 1]^d$ , where

$$71 \quad \phi(\mathbf{x}) = 2\omega_f(\sqrt{d}) \sum_{j=1}^N 2^{-j} \sigma_3 \left( a_j \sigma_2 \left( 1 + \sum_{i=1}^d 2^{(i-1)N} \sigma_1(2^N x_i) \right) \right) + f(\mathbf{0}) - \omega_f(\sqrt{d}) \quad (1)$$

72 can be implemented by an FNN with activation functions  $\sigma_1$ ,  $\sigma_2$ , and  $\sigma_3$ ,  
 73 width  $\max\{N, d\}$ , three hidden layers, and  $2(d + N + 1)$  nonzero parameters.

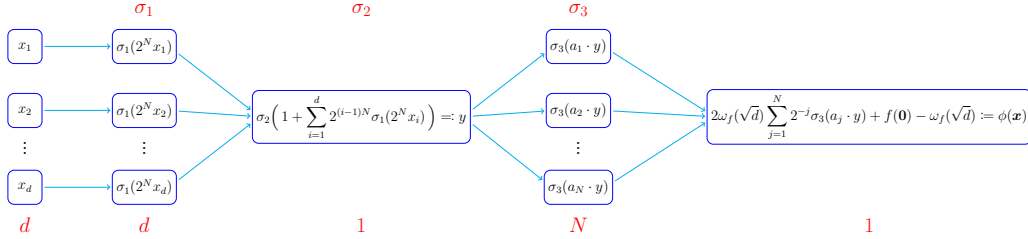


Figure 1: An illustration of the desired three-hidden-layer network in Theorem 1.1 for any  $\mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}$ . Each of the red functions “ $\sigma_1$ ”, “ $\sigma_2$ ”, and “ $\sigma_3$ ” above the network is the activation function of the corresponding hidden layer. The number of neurons in hidden layer is indicated by the red number below it.

74 In Theorem 1.1,  $[0, 1]^d$  and  $[0, 1]^d$  can be unified and generalized to  
 75  $[0, M]^d$  for any  $M > 0$  at the price of more complicated constants in the ap-  
 76 proximation rate as we shall see in a corollary later. The rate in  $\omega_f(\sqrt{d} 2^{-N})$   
 77 implicitly depends on  $N$  through the modulus of continuity of  $f$ , while the  
 78 rate in  $2\omega_f(\sqrt{d}) 2^{-N}$  is explicit in  $N$ . Simplifying the implicit approximation  
 79 rate to make it explicitly depending on  $N$  is challenging in general. How-  
 80 ever, if  $f$  is a Hölder continuous function on  $[0, 1]^d$  of order  $\alpha \in (0, 1]$  with a  
 81 constant  $\lambda$ , i.e.,  $f(\mathbf{x})$  satisfying

$$82 \quad |f(\mathbf{x}) - f(\mathbf{y})| \leq \lambda \|\mathbf{x} - \mathbf{y}\|_2^\alpha, \quad \text{for any } \mathbf{x}, \mathbf{y} \in [0, 1]^d, \quad (2)$$

83 then  $\omega_f(r) \leq \lambda r^\alpha$  for any  $r \geq 0$ . Therefore, in the case of Hölder continuous  
 84 functions, the approximation rate is simplified to  $3\lambda d^{\alpha/2} 2^{-\alpha N}$  as shown in the  
 85 following corollary. In the special case of Lipschitz continuous functions with  
 86 a Lipschitz constant  $\lambda$ , the approximation rate is simplified to  $3\lambda \sqrt{d} 2^{-N}$ .

87 **Corollary 1.2.** *Given any  $N \in \mathbb{N}^+$  and a Hölder continuous function  $f$  on*  
 88  *$[0, 1]^d$  of order  $\alpha$  with a constant  $\lambda$ , there exists a function  $\phi$  implemented*

89 by a three-hidden-layer FNN with activation functions  $\sigma_1$ ,  $\sigma_2$ , and  $\sigma_3$ , width  
 90  $\max\{d, N\}$  and  $2(d + N + 1)$  nonzero parameters such that

$$91 \quad |\phi(\mathbf{x}) - f(\mathbf{x})| \leq 3\lambda d^{\alpha/2} 2^{-\alpha N}, \quad \text{for any } \mathbf{x} \in [0, 1]^d.$$

92 First, Theorem 1.1 and its corollaries show that the approximation ca-  
 93 pacity of three-hidden-layer neural networks with simple activation functions  
 94 for continuous functions can be exponentially improved by increasing the net-  
 95 work width, and the approximation error can be explicitly characterized in  
 96 terms of the width  $\mathcal{O}(N)$ . Second, this new class of networks overcomes the  
 97 curse of dimensionality in the approximation power when the modulus of con-  
 98 tinuity is moderate, since the approximation order is essentially  $\sqrt{d}$  times a  
 99 function of  $N$  independent of  $d$  within the modulus of continuity. Therefore,  
 100 three hidden layers are enough for neural networks to achieve exponential  
 101 convergence and avoid the curse of dimensionality for generic functions. The  
 102 width is also powerful in network approximation.

103 The rest of this paper is organized as follows. In Section 2, we discuss the  
 104 application scope of our theory and compare related works in the literature.  
 105 We will prove Theorem 1.1 and its corollaries in Section 3. Finally, we  
 106 conclude this paper in Section 4.

## 107 2. Discussion

108 In this section, we will further interpret our results and discuss related  
 109 research in the field of neural network approximation.

### 110 2.1. Application scope of our theory in machine learning

111 In supervised learning, an unknown target function  $f(\mathbf{x})$  defined on a  
 112 domain  $\Omega$  is learned through its finitely many samples  $\{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n$ . If  
 113 neural networks are applied in supervised learning, the following optimization  
 114 problem is solved to identify a neural network  $\phi(\mathbf{x}; \boldsymbol{\theta}_S)$  with  $\boldsymbol{\theta}_S$  as the set of  
 115 parameters to infer  $f(\mathbf{x})$  for unseen data samples  $\mathbf{x}$ :

$$116 \quad \boldsymbol{\theta}_S = \arg \min_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}) := \arg \min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{\{\mathbf{x}_i\}_{i=1}^n} \ell(\phi(\mathbf{x}_i; \boldsymbol{\theta}), f(\mathbf{x}_i)) \quad (3)$$

117 with a loss function typically taken as  $\ell(y, y') = \frac{1}{2}|y - y'|^2$ . The inference error  
 118 is usually measured by  $R_D(\boldsymbol{\theta}_S)$ , where

$$119 \quad R_D(\boldsymbol{\theta}) := \mathbb{E}_{\mathbf{x} \sim U(\Omega)} [\ell(\phi(\mathbf{x}; \boldsymbol{\theta}), f(\mathbf{x}))],$$

where the expectation is taken with an unknown data distribution  $U(\Omega)$  over  $\Omega$ . In the analysis,  $U(\Omega)$  is assumed to be known, e.g, a uniform distribution for simplicity.

Note that the best neural network to infer  $f(\mathbf{x})$  is  $\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{D}})$  with  $\boldsymbol{\theta}_{\mathcal{D}}$  given by

$$\boldsymbol{\theta}_{\mathcal{D}} = \arg \min_{\boldsymbol{\theta}} R_{\mathcal{D}}(\boldsymbol{\theta}).$$

The best possible inference error is  $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})$ . In real applications,  $U(\Omega)$  is unknown and only finitely many samples from this distribution are available. Hence, the empirical loss  $R_{\mathcal{S}}(\boldsymbol{\theta})$  is minimized hoping to obtain  $\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{S}})$ , instead of minimizing the population loss  $R_{\mathcal{D}}(\boldsymbol{\theta})$  to obtain  $\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{D}})$ . In practice, a numerical optimization method to solve (3) may result in a numerical solution (denoted as  $\boldsymbol{\theta}_{\mathcal{N}}$ ) that may not be a global minimizer  $\boldsymbol{\theta}_{\mathcal{S}}$ . Therefore, the actually learned neural network to infer  $f(\mathbf{x})$  is  $\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{N}})$  and the corresponding inference error is measured by  $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}})$ .

By the discussion just above, it is crucial to quantify  $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}})$  to see how good the learned neural network  $\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{N}})$  is, since  $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}})$  is the expected inference error over all possible data samples. Note that

$$\begin{aligned} R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}}) &= [R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}})] + [R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}})] + [R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{D}})] \\ &\quad + [R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{D}}) - R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})] + R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}}) \\ &\leq R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}}) + [R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}})] \\ &\quad + [R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}})] + [R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{D}}) - R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})], \end{aligned} \tag{4}$$

where the inequality comes from the fact that  $[R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{D}})] \leq 0$  since  $\boldsymbol{\theta}_{\mathcal{S}}$  is a global minimizer of  $R_{\mathcal{S}}(\boldsymbol{\theta})$ . The constructive approximation established in this paper and in the literature provides an upper bound of  $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})$  in terms of the network size, e.g., in terms of the network width and depth, or in terms of the number of parameters. The second term of (4) is bounded by the optimization error of the numerical algorithm applied to solve the empirical loss minimization problem in (3). If the numerical algorithm is able to find a global minimizer, the second term is equal to zero. The theoretical guarantee of the convergence of an optimization algorithm to a global minimizer  $\boldsymbol{\theta}_{\mathcal{S}}$  and the characterization of the convergence belong to the optimization analysis of neural networks. The third and fourth term of (4) are usually bounded in terms of the sample size  $n$  and a certain norm of the corresponding set of parameters  $\boldsymbol{\theta}_{\mathcal{N}}$  and  $\boldsymbol{\theta}_{\mathcal{D}}$ , respectively. The study of the bounds for the third

and fourth terms is referred to as the generalization error analysis of neural networks.

Theorem 1.1 and its corollaries provide an upper bound of  $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})$ . This bound only depends on the given budget of neurons and layers of FLES networks. Hence, this bound is independent of the empirical loss minimization in (3) and the optimization algorithm used to compute the numerical solution of (3). In other words, Theorem 1.1 and its corollaries quantify the approximation power of FLES networks with a given size. Designing efficient optimization algorithms and analyzing the generalization bounds for FLESU networks are two other separate future directions.

## 2.2. Further interpretation of our theory

In the interpretation of our theory, two more aspects are important to discuss. The first one is whether it is possible to extend our theory to functions on a more general domain, e.g,  $[0, M]^d$  for any  $M > 0$ , because  $M > 1$  may cause an implicit curse of dimensionality in some existing theory. The second one is how bad the modulus of continuity would be since it is related to a high-dimensional function  $f$  that may lead to an implicit curse of dimensionality in our approximation rate.

First, we can generalize Theorem 1.1 to the function space  $C([0, M]^d)$  for any  $M > 0$  in the following corollary with the modulus of continuity  $\omega_f^{[0, M]^d}$  defined as follows. For an arbitrary set  $E \subset \mathbb{R}^d$ ,  $\omega_f^E(r)$  is defined via  $\omega_f^E(r) := \sup \{ |f(\mathbf{x}) - f(\mathbf{y})| : \|\mathbf{x} - \mathbf{y}\|_2 \leq r, \mathbf{x}, \mathbf{y} \in E \}$ , for any  $r \geq 0$ . As defined earlier,  $\omega_f(r)$  is short of  $\omega_f^{[0, 1]^d}(r)$ . The proof of this corollary will be presented in Section 3.2.

**Corollary 2.1.** *Given any  $N \in \mathbb{N}^+$  and a continuous function  $f$  on  $[0, M]^d$  for any  $M > 0$ , there exist  $a_1, a_2, \dots, a_N \in [0, \frac{1}{2})$  such that*

$$|\phi(\mathbf{x}) - f(\mathbf{x})| \leq 2\omega_f^{[0, M]^d}(2M\sqrt{d})2^{-N} + \omega_f^{[0, M]^d}(2M\sqrt{d})2^{-N},$$

for any  $\mathbf{x} = (x_1, \dots, x_d) \in [0, M]^d$ , where

$$\phi(\mathbf{x}) = 2\omega_f^{[0, M]^d}(2M\sqrt{d}) \sum_{j=1}^N 2^{-j} \sigma_3 \left( a_j \sigma_2 \left( 1 + \sum_{i=1}^d 2^{(i-1)N} \sigma_1 \left( \frac{2^N x_i}{2M} \right) \right) \right) + f(\mathbf{0}) - \omega_f^{[0, M]^d}(2M\sqrt{d})$$

can be implemented by a three-hidden-layer FNN with activation functions  $\sigma_1, \sigma_2$ , and  $\sigma_3$ , width  $\max\{N, d\}$ , and  $2(d + N + 1)$  nonzero parameters.

185 Hence, the size of the function domain  $[0, M]^d$  only has a mild influence  
 186 on the approximation rate of our FLES networks. FLES networks can still  
 187 avoid the curse of dimensionality and achieve exponential convergence for  
 188 continuous functions on  $[0, M]^d$  when  $M > 1$ . For example, in the case of  
 189 Hölder continuous functions of order  $\alpha$  with a constant  $\lambda$  on  $[0, M]^d$ , our  
 190 approximation rate becomes  $3\lambda(2M\sqrt{d}2^{-N})^\alpha$ .

191 Second, most interesting continuous functions in practice have a good  
 192 modulus of continuity such that there is no implicit curse of dimensionality  
 193 hidden in  $\omega_f(\cdot)$ . For example, we have discussed the case of Hölder contin-  
 194 uous functions previously. We would like to remark that the class of Hölder  
 195 continuous functions implicitly depends on  $d$  through its definition in (2),  
 196 but this dependence is moderate since the  $\ell^2$ - norm in (2) is the square root  
 197 of a sum with  $d$  terms. Let us now discuss several cases of  $\omega_f(\cdot)$  when we  
 198 cannot achieve exponential convergence or cannot avoid the curse of dimen-  
 199 sionality. The first example is  $\omega_f(r) = \frac{1}{\ln(1/r)}$  for a small  $r > 0$ , which leads to  
 200 an approximation rate

$$201 \quad 3(N \ln 2 - \tfrac{1}{2} \ln d)^{-1}, \quad \text{for large } N \in \mathbb{N}^+.$$

202 Apparently, the above approximation rate still avoids the curse of dimension-  
 203 ality but there is no exponential convergence, which has been canceled out  
 204 by “ln” in  $\omega_f(\cdot)$ . The second example is  $\omega_f(r) = \frac{1}{\ln^{1/d}(1/r)}$  for a small  $r > 0$ ,  
 205 which leads to an approximation rate

$$206 \quad 3(N \ln 2 - \tfrac{1}{2} \ln d)^{-1/d}, \quad \text{for large } N \in \mathbb{N}^+.$$

207 The power  $\frac{1}{d}$  further weakens the approximation rate and hence the curse of  
 208 dimensionality exists. The last example we would like to discuss is  $\omega_f(r) =$   
 209  $r^{\alpha/d}$  for a small  $r > 0$ , which results in the approximation rate

$$210 \quad 3\lambda d^{\frac{\alpha}{2d}} 2^{-\frac{\alpha}{d}N}, \quad \text{for large } N \in \mathbb{N}^+,$$

211 which achieves the exponential convergence and avoids the curse of dimen-  
 212 sionality when we use very deep networks with a fixed width. But if we fix  
 213 depth, there is no exponential convergence and the curse exists. Though  
 214 we have provided several examples of immoderate  $\omega_f(\cdot)$ , to the best of our  
 215 knowledge, we are not aware of useful continuous functions with  $\omega_f(\cdot)$  that  
 216 is immoderate.



### 2.3. Kolmogorov-Arnold Superposition Theorem

A closely related research topic is the Kolmogorov-Arnold representation theorem (KST) [Kolmogorov \(1956\)](#); [Arnold \(1957\)](#); [Kolmogorov \(1957\)](#) and its approximation in a form of modern neural networks. Our FLES networks admit super approximation power with a fixed number of layers for continuous functions and the KST exactly represent continuous functions using two hidden layers and  $\mathcal{O}(d)$  neurons. More specifically, given any  $f \in C([0, 1]^d)$ , the KST shows that there exist continuous functions  $\phi_q : \mathbb{R} \rightarrow \mathbb{R}$  and  $\psi_{q,p} : \mathbb{R} \rightarrow \mathbb{R}$  such that

$$f(\mathbf{x}) = \sum_{q=0}^{2d} \phi_q \left( \sum_{p=1}^d \psi_{q,p}(x_p) \right), \quad \text{for any } \mathbf{x} = (x_1, \dots, x_d) \in [0, 1]^d. \quad (5)$$

Note that the activation functions  $\{\phi_q\}$  (also called outer functions) of the neural network in Equation (5) have to depend on the target function  $f$ , though  $\{\psi_{q,p}\}$  (also called inner functions) can be independent of  $f$ . The modulus of continuity of  $\{\psi_{q,p}\}$  can be constructed such that they moderately depend on  $d$ , but the modulus of continuity of  $\{\phi_q\}$  would be exponentially bad in  $d$ . In sum, the outer functions are too pathological such that there is no existing numerical algorithms to evaluate these activation functions, even though they are shown to exist by iterative construction [Braun and Griebel \(2009\)](#).

There has been an active research line to develop more practical network approximation based on KST [Kůrková \(1991, 1992\)](#); [Maierov and Pinkus \(1999\)](#); [Guliyev and Ismailov \(2018\)](#); [Montanelli and Yang \(2020\)](#); [Igel'nik and Parikh \(2003\)](#); [Schmidt-Hieber \(2020\)](#) by relaxing the exact representation to network approximation with an  $\varepsilon$ -error. The key issue these KST-related networks attempting to address is the  $f$ -dependency of the activation functions and the main goal is to construct neural networks conquering the curse of dimensionality in a more practical way computationally. The main idea of these variants is to apply computable activation functions independent of  $f$  to construct neural networks to approximate the outer and inner functions of the KST, resulting in a larger network that can approximate a continuous function with the desired accuracy. Using this idea, the seminal work in [Kůrková \(1992\)](#) applied sigmoid activation functions and constructed two-hidden-layer networks to approximate  $f \in C([0, 1]^d)$ . Though the activation functions are independent of  $f$ , the number of neurons scales exponentially in  $d$  and the curse of dimensionality exists. Cubic-splines and piecewise linear

functions have also been used to approximate the outer and inner functions of KST in [Igel'nik and Parikh \(2003\)](#); [Montanelli and Yang \(2020\)](#); [Schmidt-Hieber \(2020\)](#), resulting in cubic-spline networks or deep ReLU networks to approximate  $f \in C([0, 1]^d)$ . But the approximation bounds in these works still suffer from the curse of dimensionality unless  $f$  has simple outer functions in the KST. It is still an open problem to characterize the class of functions with a moderate outer function in KST.

To the best of our knowledge, the most successful construction of neural networks with  $f$ -independent activation functions conquering the curse of dimensionality is in [Maierov and Pinkus \(1999\)](#); [Guliyev and Ismailov \(2018\)](#), where a two-hidden-layer network with  $\mathcal{O}(d)$  neurons can approximate  $f \in C([0, 1]^d)$  within an arbitrary error  $\varepsilon$ . Let us briefly summarize their main ideas to obtain such an exciting result here. 1) Identify a dense and countable subset  $\{u_k\}_{k=1}^\infty$  of  $C([-1, 1])$ , e.g., polynomials with rational coefficients. 2) Construct an activation function  $\varrho$  to “store” all  $u_k(x)$  for  $x \in [-1, 1]$ . For example, divide the domain of  $\varrho(x)$  into countable pieces and each piece is a connected interval of length 2 associated with a  $u_k$ . In particular, let  $\varrho(x + 4k + 1) = a_k + b_k x + c_k u_k(x)$  for any  $x \in [-1, 1]$  with carefully chosen constants  $a_k$ ,  $b_k$ , and  $c_k$  such that  $\varrho(x)$  can be a sigmoid function. 3) By construction, there exists a one-hidden layer network with width 3 and  $\varrho(x)$  as the activation function to approximate any outer or inner function in KST with an arbitrary accuracy parameter  $\delta$ . Only the parameters of the one-hidden-layer network depend on the target function and accuracy. 4) Replace the inner and outer function in KST with these one-hidden-layer networks to achieve a two-hidden-layer network with  $\varrho(x)$  as the activation function and width  $\mathcal{O}(d)$  to approximate an arbitrary  $f \in C([0, 1]^d)$  within an arbitrary error  $\varepsilon$ . Unfortunately, the construction of the parameters of this magic network relies on the evaluation of the outer and inner functions of KST, which is not computationally feasible even if computation with arbitrary precision is allowed.

We would like to remark that, though the approximation rate of FLES networks in this paper is relatively worse than the approximation rate in [Maierov and Pinkus \(1999\)](#); [Guliyev and Ismailov \(2018\)](#), our activation functions are much simpler and there are explicit formulas to specify the parameters of FLES networks. If computation with an arbitrary precision is allowed and the target function  $f$  can be arbitrarily sampled, we can specify all the weights in FLES networks. Besides, our approximation rate is sufficiently attractive since it is exponential and avoids the curse of dimen-

Table 1: A comparison of several KST-related results for approximating  $f \in C([0, 1]^d)$ .

paper	number of hidden layers	width	activation function(s)	error	remark
Kolmogorov (1956); Arnold (1957); Kolmogorov (1957)	2	$2d + 1$	$f$ -dependent	0	original KST
Majorov and Pinkus (1999); Guliyev and Ismailov (2018)	2	$\mathcal{O}(d)$	$f$ -independent	arbitrary error $\varepsilon$	based on KST
Shen et al. (2019)	3	$\mathcal{O}(dN)$	ReLU	$\mathcal{O}(N^{-2/d})$	not based on KST
this paper	3	$\max\{d, N\}$	$(\sigma_1, \sigma_2, \sigma_3)$	$2\omega_f(\sqrt{d})2^{-N} + \omega_f(\sqrt{d})2^{-N}$	not based on KST

sionality. For a large dimension  $d$ , the width parameter of our FLES network can be chosen as  $N = d$ , which leads to a FLES network of size  $\mathcal{O}(d)$  with an approximation accuracy  $\mathcal{O}(2^{-d})$  for Lipschitz continuous functions.  $\mathcal{O}(2^{-d})$  is sufficiently attractive. In practice, when  $d$  is very large,  $N$  could be much smaller than  $d$  and our approximation rate is still attractive.

Finally, we list several KST-related results in Table 1 for a quick comparison.

#### 2.4. Discussion on the literature

In this section, we will discuss other recent development of neural network approximation. Our discussion will be divided into mainly three parts according to the analysis methodology in the references: 1) functions admitting integral representations; 2) linear approximation; 3) bit extraction.

In the seminal work of Barron (1993), its variants or generalization Barron and Klusowski (2018); E et al. (2019); Chen and Wu (2019); Montanelli et al. (2020), and related references therein,  $d$ -dimensional functions of the following form were considered:

$$f(\mathbf{x}) = \int_{\tilde{\Omega}} a(\mathbf{w}) K(\mathbf{w} \cdot \mathbf{x}) d\mu(\mathbf{w}), \quad (6)$$

where  $\tilde{\Omega} \subset \mathbb{R}^d$ ,  $\mu(\mathbf{w})$  is a Lebesgue measure in  $\mathbf{w}$ , and  $\mathbf{x} \in \Omega \subset \mathbb{R}^d$ . The above integral representation is equivalent to the expectation of a high-dimensional random function when  $\mathbf{w}$  is treated as a random variable. By the law of large number theory, the average of  $N$  samples of the integrand leads to an approximation of  $f(\mathbf{x})$  with an approximation error bounded by  $\frac{C_f \sqrt{\mu(\Omega)}}{\sqrt{N}}$  measured in  $L^2(\Omega, \mu)$  (Equation (6) of Barron (1993)), where  $\mathcal{O}(N)$  is the total number of parameters in the network,  $C_f$  is a  $d$ -dimensional integral with an integrand related to  $f$ , and  $\mu(\Omega)$  is the Lebesgue measure of  $\Omega$ . As discussed in Barron (1993),  $\mu(\Omega)$  and  $C_f$  would be exponential in  $d$  and standard smoothness properties of  $f$  alone are not enough to remove the exponential dependence of  $C_f$  on  $d$ . Therefore, the curse of dimensionality

exists in the whole approximation error while the curse does not exist in the approximation rate in  $N$ .

Linear approximation is an efficient approximation tool for smooth functions that computes the approximant of a target function via a linear projection to a Hilbert space or a Banach space as the approximant space. Typical examples include approximation via orthogonal polynomials, Fourier series expansion, etc. Inspired by the seminal work in [Yarotsky \(2017\)](#), where deep ReLU networks were constructed to approximate polynomials with exponential convergence, subsequent works in [E and Wang \(2018\)](#); [Opschoor et al. \(2019\)](#); [Montanelli and Du \(2019\)](#); [Chen and Wu \(2019\)](#); [Montanelli et al. \(2020\)](#); [Yarotsky and Zhevnerchuk \(2019\)](#); [Lu et al. \(2020\)](#); [Montanelli and Yang \(2020\)](#); [Yang and Wang \(2020\)](#) have constructed deep ReLU networks to approximate various smooth function spaces. The main idea of these works is to approximate smooth functions via (piecewise) polynomial approximation first and then construct deep ReLU networks to approximate the ensemble of polynomials. If the approximation rate of polynomials to approximate the target function is exponential, then the approximation rate of deep ReLU networks to approximate the target function is also exponential. But the curse of dimensionality exists since polynomial approximation cannot avoid the curse.

The bit extraction proposed in [Bartlett et al. \(1998\)](#) has been a very important technique to develop nearly optimal approximation rates of deep ReLU neural networks [Yarotsky \(2018\)](#); [Shen et al. \(2020\)](#); [Lu et al. \(2020\)](#); [Yang and Wang \(2020\)](#) and the optimality is based on the nearly optimal VC-dimension bound of ReLU networks in [Harvey et al. \(2017\)](#). The bit extraction was also applied in [Shen et al. \(2020\)](#); [Schmidt-Hieber \(2020\)](#) and this paper to develop network approximation theories. In the first step, an efficient projection map in a form of a ReLU, or a Floor-ReLU, or a FLES network is constructed to project high-dimensional points to one-dimensional points such that the high-dimensional approximation problem is reduced to a one-dimensional approximation problem. In the first step, the one-dimensional approximation problem is solved by constructing a ReLU, or a Floor-ReLU, or a FLES network, which can be efficiently compressed via the bit extraction. Although shallower neural networks can also carry out the above two steps, bit extraction can take full advantage of the power of depth and construct deep neural networks with a nearly optimal number of parameters or neurons to fulfill the above two steps.

### 355 3. Theoretical Analysis

356 In this section, we first introduce basic notations in this paper in Section  
 357 3.1. Then we prove Theorem 1.1 and its corollaries in Section 3.2.

#### 358 3.1. Notations

359 The main notations of this paper are listed as follows.

- 360 • Vectors and matrices are denoted in a bold font. Standard vectorization  
 361 is adopted in the matrix and vector computation. For example, a scalar  
 362 plus a vector means adding the scalar to each entry of the vector.
- 363 • Let  $\mathbb{N}^+$  denote the set containing all positive integers, i.e.,  $\mathbb{N}^+ = \{1, 2, 3, \dots\}$ .
- 364 • For any  $p \in [1, \infty)$ , the  $p$ -norm of a vector  $\mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$  is  
 365 defined by

$$366 \quad \|\mathbf{x}\|_p := (|x_1|^p + |x_2|^p + \dots + |x_d|^p)^{1/p}.$$

- 367 • Let  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  denote the rectified linear unit (ReLU), i.e.  $\sigma(x) =$   
 368  $\max\{0, x\}$ . With a slight abuse of notation, we define  $\sigma : \mathbb{R}^d \rightarrow \mathbb{R}^d$  as

$$369 \quad \sigma(\mathbf{x}) = \begin{bmatrix} \max\{0, x_1\} \\ \vdots \\ \max\{0, x_d\} \end{bmatrix} \text{ for any } \mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d.$$

- 370 • The floor function (Floor) is defined as  $\lfloor x \rfloor := \max\{n : n \leq x, n \in \mathbb{Z}\}$  for  
 371 any  $x \in \mathbb{R}$ .
- 372 • For  $\theta \in [0, 1)$ , suppose its binary representation is  $\theta = \sum_{\ell=1}^{\infty} \theta_{\ell} 2^{-\ell}$  with  
 373  $\theta_{\ell} \in \{0, 1\}$ , we introduce a special notation  $\text{bin}0.\theta_1\theta_2\cdots\theta_L$  to denote the  
 374  $L$ -term binary representation of  $\theta$ , i.e.,  $\text{bin}0.\theta_1\theta_2\cdots\theta_L := \sum_{\ell=1}^L \theta_{\ell} 2^{-\ell}$ .
- 375 • The expression “a network with a width  $N$  and a depth  $L$ ” means
  - 376 – The maximum width of this network for all **hidden** layers is no  
 377 more than  $N$ .
  - 378 – The number of **hidden** layers of this network is no more than  $L$ .

379 *3.2. Proof of Theorem 1.1 and Corollary 2.1*

380 To prove Theorem 1.1, we first present the proof sketch. Shortly speak-  
 381 ing, we construct piecewise constant functions to approximate continuous  
 382 functions. There are five key steps in our construction.

- 383 1. Normalize  $f$  as  $\tilde{f}$  satisfying  $\tilde{f}(\mathbf{x}) \in [0, 1]$  for any  $\mathbf{x} \in [0, 1]^d$ , divide  
 384  $[0, 1]^d$  into a set of non-overlapping cubes  $\{Q_\beta\}_{\beta \in \{0, 1, \dots, J-1\}^d}$ , and denote  
 385  $\mathbf{x}_\beta$  as the vertex of  $Q_\beta$  with minimum  $\|\cdot\|_1$  norm, where  $J$  is an integer  
 386 determined later. See Figure 2 for the illustrations of  $Q_\beta$  and  $\mathbf{x}_\beta$ .
- 387 2. Construct a vector-valued function  $\Phi_1 : \mathbb{R}^d \rightarrow \mathbb{R}^d$  projecting the whole  
 388 cube  $Q_\beta$  to the index  $\beta$  for each  $\beta \in \{0, 1, \dots, J-1\}^d$ , i.e.,  $\Phi_1(\mathbf{x}) = \beta$   
 389 for all  $\mathbf{x} \in Q_\beta$ .
- 390 3. Construct a linear function  $\phi_2 : \mathbb{R}^d \rightarrow \mathbb{R}$  bijectively mapping  $\beta \in$   
 391  $\{0, 1, \dots, J-1\}^d$  to  $\phi_2(\beta) \in \{1, 2, \dots, J^d\}$ .
- 392 4. Construct a function  $\phi_3 : \mathbb{R} \rightarrow \mathbb{R}$  mapping  $\phi_2(\beta) \in \{1, 2, \dots, J^d\}$  approx-  
 393 imately to  $\tilde{f}(\mathbf{x}_\beta)$ , i.e.,  $\phi_3(\phi_2(\beta)) \approx \tilde{f}(\mathbf{x}_\beta)$  for each  $\beta \in \{0, 1, \dots, J-1\}^d$ .
- 394 5. Define  $\tilde{\phi} := \phi_3 \circ \phi_2 \circ \Phi_1$ . Then  $\tilde{\phi}$  is a piecewise constant function mapping  
 395  $\mathbf{x} \in Q_\beta$  to  $\phi_3(\phi_2(\beta)) \approx \tilde{f}(\mathbf{x}_\beta)$  for each  $\beta \in \{0, 1, \dots, J-1\}^d$ , implying  
 396  $\tilde{\phi} \approx \tilde{f}$ . Finally, re-scale and shift  $\tilde{\phi}$  to obtain the final function  $\phi$   
 397 approximating  $f$  well.

398 The most technical step above is Step 4, the realization of which relies  
 399 on the proposition below.

400 **Proposition 3.1.** *Given any  $K \in \mathbb{N}^+$  and arbitrary  $\theta_1, \theta_2, \dots, \theta_K \in \{0, 1\}$ , it*  
 401 *holds that*

402 
$$\sigma_3(2^k \cdot a) = \theta_k, \quad \text{for any } k \in \{1, 2, \dots, K\},$$

403 *where*

404 
$$a = \sum_{j=1}^K 2^{-j-1} \cdot \theta_j \in [0, \frac{1}{2}).$$

405 *Proof of Proposition 3.1.* Since  $\theta_j \in \{0, 1\}$  for  $j \in \{1, 2, \dots, K\}$ , we have

406 
$$0 \leq \sum_{j=1}^K 2^{-j-1} \cdot \theta_j \leq \sum_{j=1}^K 2^{-j-1} < \frac{1}{2},$$

407 implying  $a \in [0, \frac{1}{2})$ .

408 Next, fix  $k \in \{1, 2, \dots, K\}$  for the proof below. It holds that

$$409 \quad 2^k \cdot a = 2^k \cdot \sum_{j=1}^K 2^{-j-1} \cdot \theta_j = \underbrace{\sum_{j=1}^{k-1} 2^{k-j-1} \cdot \theta_j}_{\text{an integer}} + \overbrace{\frac{1}{2} \theta_k}^{0 \text{ or } \frac{1}{2}} + \underbrace{\sum_{j=k+1}^K 2^{k-j-1} \cdot \theta_j}_{\text{in } [0, \frac{1}{2})}. \quad (7)$$

410 Clearly, the first term in Equation (7)  $\sum_{j=1}^{k-1} 2^{k-j-1} \cdot \theta_j$  is a non-negative integer  
 411 since  $\theta_j \in \{0, 1\}$  for any  $j \in \{1, 2, \dots, K\}$ . As for the third term in Equation (7),  
 412 we have

$$413 \quad 0 \leq \sum_{j=k+1}^K 2^{k-j-1} \cdot \theta_j \leq \sum_{j=k+1}^K 2^{k-j-1} < \frac{1}{2}$$

414 Therefore, By Equation (7), we have

$$415 \quad 2^k \cdot a \in \bigcup_{n \in \mathbb{N}} [n, n + \frac{1}{2}), \text{ if } \theta_k = 0, \quad \text{and} \quad 2^k \cdot a \in \bigcup_{n \in \mathbb{N}} [n + \frac{1}{2}, n + 1), \text{ if } \theta_k = 1. \quad (8)$$

416 Recall that  $\sigma_3(x) = \mathcal{T}(x - \lfloor x \rfloor - \frac{1}{2})$ , where  $\mathcal{T}(x) = \begin{cases} 0, & x \geq 0, \\ 1, & x < 0. \end{cases}$ . It is easy to verify  
 417 that

$$418 \quad \sigma_3(x) = 0 \text{ if } x \in \bigcup_{n \in \mathbb{N}} [n, n + \frac{1}{2}), \quad \text{and} \quad \sigma_3(x) = 1 \text{ if } x \in \bigcup_{n \in \mathbb{N}} [n + \frac{1}{2}, n + 1).$$

419 If  $\theta_k = 0$ , by Equation (8), we have

$$420 \quad 2^k \cdot a \in \bigcup_{n \in \mathbb{N}} [n, n + \frac{1}{2}) \implies \sigma_3(2^k \cdot a) = 0 = \theta_k.$$

421 Similarly, if  $\theta_k = 1$ , by Equation (8), we have

$$422 \quad 2^k \cdot a \in \bigcup_{n \in \mathbb{N}} [n + \frac{1}{2}, n + 1) \implies \sigma_3(2^k \cdot a) = 1 = \theta_k.$$

423 Since  $k \in \{1, 2, \dots, K\}$  is arbitrary, we have  $\sigma_3(2^k \cdot a) = \theta_k$  for any  $k \in$   
 424  $\{1, 2, \dots, K\}$ . So we finish the proof.  $\square$

---

<sup>1</sup>By convention,  $\sum_{j=n}^m a_j = 0$  if  $n > m$  no matter what  $a_j$  is for each  $j$ .

425 We would like to point out that Proposition 3.1 indicates that the VC-  
 426 dimension of the function space

$$427 \quad \{f : f(x) = \sigma_3(a \cdot x), \text{ for } a \in \mathbb{R}\}$$

428 is infinity.

429 With Proposition 3.1 in hand, we are ready to prove Theorem 1.1.

430 *Proof of Theorem 1.1.* The proof consists of five steps.

431 **Step 1:** Set up.

432 Assume  $f$  is not a constant function since it is a trivial case. Then  
 433  $\omega_f(r) > 0$  for any  $r > 0$ . Clearly,  $|f(\mathbf{x}) - f(\mathbf{0})| \leq \omega_f(\sqrt{d})$  for any  $\mathbf{x} \in [0, 1]^d$ .  
 434 Define

$$435 \quad \tilde{f} := (f - f(\mathbf{0}) + \omega_f(\sqrt{d})) / (2\omega_f(\sqrt{d})). \quad (9)$$

436 It follows that  $\tilde{f}(\mathbf{x}) \in [0, 1]$  for any  $\mathbf{x} \in [0, 1]^d$ .

437 Set  $J = 2^N$  and divide  $[0, 1]^d$  into  $J^d$  cubes  $\{Q_\beta\}_\beta$ . To be exact, defined  
 438  $\mathbf{x}_\beta := \beta/J$  and

$$439 \quad Q_\beta := \{\mathbf{x} = (x_1, x_2, \dots, x_d) : x_i \in [\frac{\beta_i}{J}, \frac{\beta_i+1}{J}) \text{ for } i = 1, 2, \dots, d\},$$

440 for each  $\beta = (\beta_1, \beta_2, \dots, \beta_d) \in \{0, 1, \dots, J-1\}^d$ . See Figure 2 for illustrations.

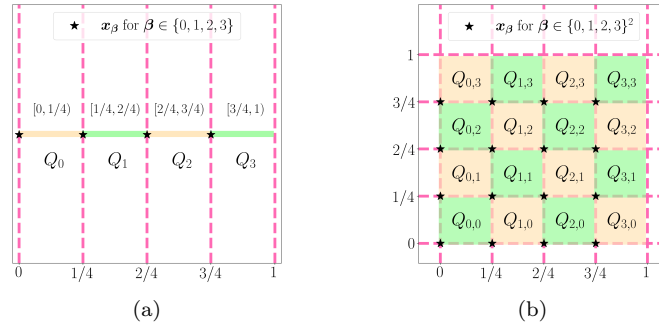


Figure 2: Illustrations of  $Q_\beta$  and  $\mathbf{x}_\beta$  for  $\beta \in \{0, 1, \dots, J-1\}^d$ . (a)  $J = 4$ ,  $d = 1$ . (b)  $J = 4$ ,  $d = 2$ .

441 **Step 2:** Construct  $\Phi_1$  mapping  $\mathbf{x} \in Q_\beta$  to  $\beta$  for each  $\beta \in \{0, 1, \dots, J-1\}^d$ .



442 Define

$$443 \quad \Phi_1(\mathbf{x}) = \left( \lfloor Jx_1 \rfloor, \lfloor Jx_2 \rfloor, \dots, \lfloor Jx_d \rfloor \right), \quad \text{for any } \mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d.$$

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

$$445 \quad \Phi_1(\mathbf{x}) = \left( \lfloor Jx_1 \rfloor, \lfloor Jx_2 \rfloor, \dots, \lfloor Jx_d \rfloor \right) = (\beta_1, \beta_2, \dots, \beta_d) = \beta. \quad (10)$$

446 **Step 3:** Construct  $\phi_2$  bijectively mapping  $\beta \in \{0, 1, \dots, J-1\}^d$  to  $\phi_2(\beta) \in$   
 447  $\{1, 2, \dots, J^d\}$ .

448 Inspired by the  $J$ -ary representation, we define a linear map

$$449 \quad \phi_2(\mathbf{x}) := 1 + \sum_{i=1}^d J^{i-1} x_i, \quad \text{for each } \mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d.$$

450 Then  $\phi_2$  is a bijection from  $\{0, 1, \dots, J-1\}^d$  to  $\{1, 2, \dots, J^d\}$ .

451 **Step 4:** Construct  $\phi_3$  mapping  $\phi_2(\beta) \in \{1, 2, \dots, J^d\}$  approximately to  $\tilde{f}(\mathbf{x}_\beta)$ .

452

453 For each  $k \in \{1, 2, \dots, J^d\}$ , there exists a unique  $\beta \in \{0, 1, \dots, J-1\}^d$  such  
 454 that  $\phi_2(\beta) = k$ . Thus, define

$$455 \quad \xi_k := \tilde{f}(\mathbf{x}_\beta) \in [0, 1], \quad \text{for any } k \in \{1, 2, \dots, J^d\} \text{ with } k = \phi_2(\beta). \quad (11)$$

456 For each  $k \in \{1, 2, \dots, J^d\}$ , there exist  $\theta_{k,1}, \theta_{k,2}, \dots, \theta_{k,N} \in \{0, 1\}$  such that

$$457 \quad |\xi_k - \text{bin } 0.\theta_{k,1}\theta_{k,2}\dots\theta_{k,N}| \leq 2^{-N}. \quad (12)$$

458 For each  $j \in \{1, 2, \dots, N\}$ , by Proposition 3.1 (set  $K = J^d$  therein), there  
 459 exists  $a_j \in [0, \frac{1}{2})$  such that

$$460 \quad \sigma_3(2^k \cdot a_j) = \theta_{k,j}, \quad \text{for any } k \in \{1, 2, \dots, J^d\}.$$

461 Define

$$462 \quad \phi_3(x) := \sum_{j=1}^N 2^{-j} \sigma_3(a_j \sigma_2(x)) = \sum_{j=1}^N 2^{-j} \sigma_3(2^x \cdot a_j), \quad \text{for any } x \in \mathbb{R}.$$

463 We get

$$464 \quad \phi_3(k) = \sum_{j=1}^N 2^{-j} \sigma_3(2^k \cdot a_j) = \sum_{j=1}^N 2^{-j} \cdot \theta_{k,j} = \text{bin } 0.\theta_{k,1}\theta_{k,2}\dots\theta_{k,N}. \quad (13)$$

465 **Step 5:** Define  $\tilde{\phi} := \phi_3 \circ \phi_2 \circ \Phi_1$  approximating  $\tilde{f}$  well, and re-scale and  
 466 shift  $\tilde{\phi}$  to obtain  $\phi$  approximating  $f$  well.

467 Define  $\tilde{\phi} := \phi_3 \circ \phi_2 \circ \Phi_1$ , by Equation (10), (11), (12), and (13), we have,  
 468 for any  $\mathbf{x} \in Q_\beta$  and each  $\beta \in \{0, 1, \dots, J-1\}^d$  with  $k = \phi_2(\beta)$ ,

$$\begin{aligned} |\tilde{\phi}(\mathbf{x}) - \tilde{f}(\mathbf{x})| &= |\phi_3 \circ \phi_2 \circ \Phi_1(\mathbf{x}) - \tilde{f}(\mathbf{x}_\beta)| + |\tilde{f}(\mathbf{x}_\beta) - \tilde{f}(\mathbf{x})| \\ &\leq |\phi_3 \circ \phi_2(\beta) - \tilde{f}(\mathbf{x}_\beta)| + \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{J}\right) \\ &\leq |\phi_3(k) - \xi_k| + \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{J}\right) \\ &\leq |\text{bin } 0.\theta_{k,1}\theta_{k,2}\dots\theta_{k,N} - \xi_k| + \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{J}\right) \leq 2^{-N} + \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{J}\right). \end{aligned}$$

470 Finally, define  $\phi := 2\omega_f(\sqrt{d})\tilde{\phi} + f(\mathbf{0}) - \omega_f(\sqrt{d})$ . Equation (9) implies  
 471  $\omega_f(r) = 2\omega_f(\sqrt{d})\omega_{\tilde{f}}(r)$  for any  $r \geq 0$ , deducing

$$\begin{aligned} |\phi(\mathbf{x}) - f(\mathbf{x})| &= 2\omega_f(\sqrt{d})|\tilde{\phi}(\mathbf{x}) - \tilde{f}(\mathbf{x})| \\ &\leq 2\omega_f(\sqrt{d})(2^{-N} + \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{J}\right)) \\ &= 2\omega_f(\sqrt{d})2^{-N} + \omega_f\left(\frac{\sqrt{d}}{J}\right) \\ &= 2\omega_f(\sqrt{d})2^{-N} + \omega_f(\sqrt{d}2^{-N}), \end{aligned}$$

473 for any  $\mathbf{x} \in \bigcup_{\beta \in \{0,1,\dots,J-1\}^d} Q_\beta = [0, 1]^d$ . It follows from  $J = 2^N$  and the defini-  
 474 tions of  $\Phi_1$ ,  $\phi_2$ , and  $\phi_3$  that

$$\begin{aligned} \phi(\mathbf{x}) &= 2\omega_f(\sqrt{d})\phi_3 \circ \phi_2 \circ \Phi_1(\mathbf{x}) + f(\mathbf{0}) - \omega_f(\sqrt{d}) \\ &= 2\omega_f(\sqrt{d})\phi_3\left(1 + \sum_{i=1}^d J^{i-1}\sigma_1(Jx_i)\right) + f(\mathbf{0}) - \omega_f(\sqrt{d}) \\ &= 2\omega_f(\sqrt{d})\sum_{j=1}^N 2^{-j}\sigma_3\left(a_j\sigma_2\left(1 + \sum_{i=1}^d 2^{(i-1)N}\sigma_1(2^N x_i)\right)\right) + f(\mathbf{0}) - \omega_f(\sqrt{d}). \end{aligned}$$

476 So we finish the proof. □

477 Finally, we present the proof of Corollary 2.1 below.

478 *Proof of Corollary 2.1.* Given any  $f \in C([0, M]^d)$ , by Lemma 4.2 of Shen  
 479 et al. (2020) (set  $E = [0, M]^d$  and  $S = [0, 2M)^d$  therein), there exists  $g \in$   
 480  $C([0, 2M)^d)$  such that

- 481 •  $f(\mathbf{x}) = g(\mathbf{x})$  for any  $\mathbf{x} \in [0, M]^d$ ;

482 •  $\omega_f^{[0,M]^d}(r) = \omega_g^{[0,2M]^d}(r)$  for any  $r \geq 0$ .

483 Set  $\tilde{g}(\mathbf{x}) = g(2M\mathbf{x})$  for any  $\mathbf{x} \in [0,1]^d$ , then  $\tilde{g} \in C([0,1]^d)$ . By applying  
484 Theorem 1.1 to  $\tilde{g}$ , there exist  $a_1, a_2, \dots, a_N \in [0, \frac{1}{2})$  such that

485  $|\phi(2M\mathbf{x}) - g(2M\mathbf{x})| = |\tilde{\phi}(\mathbf{x}) - \tilde{g}(\mathbf{x})| \leq 2\omega_{\tilde{g}}^{[0,1]^d}(\sqrt{d})2^{-N} + \omega_{\tilde{g}}^{[0,1]^d}(\sqrt{d}2^{-N}), \quad (14)$

486 for any  $\mathbf{x} \in [0,1]^d$ , where  $\phi(2M\mathbf{x}) = \tilde{\phi}(\mathbf{x})$  and

487 
$$\tilde{\phi}(\mathbf{x}) = 2\omega_{\tilde{g}}(\sqrt{d}) \sum_{j=1}^N 2^{-j} \sigma_3 \left( a_j \sigma_2 \left( 1 + \sum_{i=1}^d 2^{(i-1)N} \sigma_1(2^N x_i) \right) \right) + \tilde{g}(\mathbf{0}) - \omega_{\tilde{g}}(\sqrt{d}).$$

488 Note that  $\omega_{\tilde{g}}^{[0,1]^d}(r) = \omega_g^{[0,2M]^d}(2Mr) = \omega_f^{[0,M]^d}(2Mr)$  for any  $r \geq 0$ . Then by  
489 Equation (14), for any  $\mathbf{x} \in [0, M]^d \subseteq [0, 2M]^d$ , we have

$$\begin{aligned} |\phi(\mathbf{x}) - f(\mathbf{x})| &= |\phi(\mathbf{x}) - g(\mathbf{x})| = \left| \tilde{\phi}\left(\frac{\mathbf{x}}{2M}\right) - \tilde{g}\left(\frac{\mathbf{x}}{2M}\right) \right| \\ &\leq 2\omega_{\tilde{g}}^{[0,1]^d}(\sqrt{d})2^{-N} + \omega_{\tilde{g}}^{[0,1]^d}(\sqrt{d}2^{-N}) \\ &= 2\omega_f^{[0,M]^d}(2M\sqrt{d})2^{-N} + \omega_f^{[0,M]^d}(2M\sqrt{d}2^{-N}), \end{aligned}$$

490

491 where

$$\begin{aligned} \phi(\mathbf{x}) &= 2\omega_f^{[0,M]^d}(2M\sqrt{d}) \sum_{j=1}^N 2^{-j} \sigma_3 \left( a_j \sigma_2 \left( 1 + \sum_{i=1}^d 2^{(i-1)N} \sigma_1\left(\frac{2^N x_i}{2M}\right) \right) \right) \\ &\quad + f(\mathbf{0}) - \omega_f^{[0,M]^d}(2M\sqrt{d}). \end{aligned}$$

492

493 With the discussion above, we have proved Corollary 2.1.  $\square$

## 494 4. Conclusion

495 This paper has introduced a theoretical framework to show that three  
496 hidden layers are enough for neural network approximation to achieve expo-  
497 nential convergence and avoid the curse of dimensionality for approximat-  
498 ing functions as general as (Hölder) continuous functions. The key idea  
499 is to leverage the power of multiple simple activation functions: the Floor  
500 function ( $\lfloor x \rfloor$ ), the exponential function ( $2^x$ ), the step function ( $\mathbb{1}_{x \geq 0}$ ), or  
501 their compositions. This new class of networks is called the FLES network.  
502 Given a Lipschitz continuous function  $f$  on  $[0,1]^d$ , it was shown by construc-  
503 tion that FLES networks with width  $\max\{d, N\}$  and three hidden layers

admit a uniform approximation rate  $3\lambda\sqrt{d}2^{-N}$ , where  $\lambda$  is the Lipschitz constant of  $f$ . More generally for an arbitrary continuous function  $f$  on  $[0, 1]^d$  with a modulus of continuity  $\omega_f(\cdot)$ , the constructive approximation rate is  $\omega_f(\sqrt{d}2^{-N}) + 2\omega_f(\sqrt{d})2^{-N}$ . The results in this paper provide a theoretical lower bound of the power of FLES networks. Whether or not this bound is achievable in actual computation relies on advanced algorithm design as a separate line of research.

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

## References

- Arnold, V.I., 1957. On functions of three variables. Dokl. Akad. Nauk SSSR, 679–681.
- Barron, A.R., 1993. Universal approximation bounds for superpositions of a sigmoidal function. IEEE Transactions on Information Theory 39, 930–945. doi:[10.1109/18.256500](https://doi.org/10.1109/18.256500).
- Barron, A.R., Klusowski, J.M., 2018. Approximation and estimation for high-dimensional deep learning networks. [arXiv:1809.03090](https://arxiv.org/abs/1809.03090).
- Bartlett, P., Maierov, V., Meir, R., 1998. Almost linear VC-dimension bounds for piecewise polynomial networks. Neural Computation 10, 217–3.
- Braun, J., Griebel, M., 2009. On a constructive proof of kolmogorov’s superposition theorem. Constructive Approximation 30, 653–675.
- Chen, L., Wu, C., 2019. A note on the expressive power of deep rectified linear unit networks in high-dimensional spaces. Mathematical Methods in the Applied Sciences 42, 3400–3404. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/mma.5575>, doi:[10.1002/mma.5575](https://doi.org/10.1002/mma.5575), [arXiv:https://onlinelibrary.wiley.com/doi/pdf/10.1002/mma.5575](https://arxiv.org/abs/https://onlinelibrary.wiley.com/doi/pdf/10.1002/mma.5575).
- Chen, M., Jiang, H., Liao, W., Zhao, T., 2019a. Efficient approximation of deep ReLU networks for functions on low dimensional manifolds, in: Wallach, H., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E.,

- 534 Garnett, R. (Eds.), Advances in Neural Information Processing Systems  
535 32. Curran Associates, Inc., pp. 8174–8184.
- 536 Chen, Z., Cao, Y., Zou, D., Gu, Q., 2019b. How much over-parameterization  
537 is sufficient to learn deep ReLU networks? CoRR arXiv:1911.12360. URL:  
538 <https://arxiv.org/abs/1911.12360>.
- 539 Du, S.S., Zhai, X., Póczos, B., Singh, A., 2019. Gradient descent prov-  
540 ably optimizes over-parameterized neural networks, in: International Con-  
541 ference on Learning Representations. URL: [https://openreview.net/](https://openreview.net/forum?id=S1eK3i09YQ)  
542 [forum?id=S1eK3i09YQ](https://openreview.net/forum?id=S1eK3i09YQ).
- 543 E, W., Ma, C., Wu, L., 2019. A priori estimates of the population risk for  
544 two-layer neural networks. Communications in Mathematical Sciences 17,  
545 1407 – 1425.
- 546 E, W., Wang, Q., 2018. Exponential convergence of the deep neural network  
547 approximation for analytic functions. CoRR abs/1807.00297. URL: [http://](http://arxiv.org/abs/1807.00297)  
548 [arxiv.org/abs/1807.00297](http://arxiv.org/abs/1807.00297), [arXiv:1807.00297](https://arxiv.org/abs/1807.00297).
- 549 Gühring, I., Kutyniok, G., Petersen, P., 2019. Error bounds for approxima-  
550 tions with deep relu neural networks in  $w^{s,p}$  norms. [arXiv:1902.07896](https://arxiv.org/abs/1902.07896).
- 551 Guliyev, N.J., Ismailov, V.E., 2018. Approximation capability of two  
552 hidden layer feedforward neural networks with fixed weights. Neu-  
553 rocomputing 316, 262 – 269. URL: [http://www.sciencedirect.](http://www.sciencedirect.com/science/article/pii/S0925231218309111)  
554 [com/science/article/pii/S0925231218309111](http://www.sciencedirect.com/science/article/pii/S0925231218309111), doi:[https://doi.org/](https://doi.org/10.1016/j.neucom.2018.07.075)  
555 [10.1016/j.neucom.2018.07.075](https://doi.org/10.1016/j.neucom.2018.07.075).
- 556 Harvey, N., Liaw, C., Mehrabian, A., 2017. Nearly-tight VC-dimension  
557 bounds for piecewise linear neural networks, in: Kale, S., Shamir, O.  
558 (Eds.), Proceedings of the 2017 Conference on Learning Theory, PMLR,  
559 Amsterdam, Netherlands. pp. 1064–1068. URL: [http://proceedings.](http://proceedings.mlr.press/v65/harvey17a.html)  
560 [mlr.press/v65/harvey17a.html](http://proceedings.mlr.press/v65/harvey17a.html).
- 561 Huttenlocher, M., Jentzen, A., Wursterberger, v.W., 2020. Overcoming the  
562 curse of dimensionality in the approximative pricing of financial derivatives  
563 with default risks. Electron. J. Probab. 25, 73 pp. URL: [https://doi.](https://doi.org/10.1214/20-EJP423)  
564 [org/10.1214/20-EJP423](https://doi.org/10.1214/20-EJP423), doi:[10.1214/20-EJP423](https://doi.org/10.1214/20-EJP423).

- 565 Igelnik, B., Parikh, N., 2003. Kolmogorov’s spline network. IEEE Transac-  
566 tions on Neural Networks 14, 725–733.
- 567 Jacot, A., Gabriel, F., Hongler, C., 2018. Neural tangent kernel: Conver-  
568 gence and generalization in neural networks, in: Bengio, S., Wallach,  
569 H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (Eds.),  
570 Advances in Neural Information Processing Systems 31. Curran As-  
571 sociates, Inc., pp. 8571–8580. URL: [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/8076-neural-tangent-kernel-convergence-and-generalization-in-neural-networks.pdf)  
572 [8076-neural-tangent-kernel-convergence-and-generalization-in-neural-networks.](http://papers.nips.cc/paper/8076-neural-tangent-kernel-convergence-and-generalization-in-neural-networks.pdf)  
573 [pdf](http://papers.nips.cc/paper/8076-neural-tangent-kernel-convergence-and-generalization-in-neural-networks.pdf).
- 574 Kolmogorov, A.N., 1956. On the representation of continuous functions of  
575 several variables by superposition of continuous functions of a smaller num-  
576 ber of variables. Dokl. Akad. Nauk SSSR , 179–182.
- 577 Kolmogorov, A.N., 1957. On the representation of continuous functions of  
578 several variables by superposition of continuous functions of one variable  
579 and addition. Dokl. Akad. Nauk SSSR , 953–956.
- 580 Kůrková, V., 1991. Kolmogorov’s theorem is relevant. Neu-  
581 ral Computation 3, 617–622. URL: [https://doi.org/10.](https://doi.org/10.1162/neco.1991.3.4.617)  
582 [1162/neco.1991.3.4.617](https://doi.org/10.1162/neco.1991.3.4.617), doi:[10.1162/neco.1991.3.4.617](https://doi.org/10.1162/neco.1991.3.4.617),  
583 [arXiv:https://doi.org/10.1162/neco.1991.3.4.617](https://doi.org/10.1162/neco.1991.3.4.617). pMID:  
584 31167327.
- 585 Kůrková, V., 1992. Kolmogorov’s theorem and multilayer neural net-  
586 works. Neural Networks 5, 501 – 506. URL: [http://www.sciencedirect.](http://www.sciencedirect.com/science/article/pii/0893608092900128)  
587 [com/science/article/pii/0893608092900128](http://www.sciencedirect.com/science/article/pii/0893608092900128), doi:[https://doi.org/](https://doi.org/10.1016/0893-6080(92)90012-8)  
588 [10.1016/0893-6080\(92\)90012-8](https://doi.org/10.1016/0893-6080(92)90012-8).
- 589 Lu, J., Shen, Z., Yang, H., Zhang, S., 2020. Deep network approximation for  
590 smooth functions. arXiv e-prints , arXiv:2001.03040[arXiv:2001.03040](https://arxiv.org/abs/2001.03040).
- 591 Lu, Y., Ma, C., Lu, Y., Lu, J., Ying, L., 2020. A mean-field analy-  
592 sis of deep resnet and beyond: Towards provable optimization via over-  
593 parameterization from depth. CoRR abs/2003.05508. URL: [https:](https://arxiv.org/abs/2003.05508)  
594 [//arxiv.org/abs/2003.05508](https://arxiv.org/abs/2003.05508), [arXiv:2003.05508](https://arxiv.org/abs/2003.05508).
- 595 Luo, T., Yang, H., 2020. Two-layer neural networks for partial differential  
596 equations: Optimization and generalization theory. ArXiv abs/2006.15733.

- 597 Maierov, V., Pinkus, A., 1999. Lower bounds for approximation by  
598 MLP neural networks. *Neurocomputing* 25, 81 – 91. URL: [http://](http://www.sciencedirect.com/science/article/pii/S0925231298001118)  
599 [www.sciencedirect.com/science/article/pii/S0925231298001118](http://www.sciencedirect.com/science/article/pii/S0925231298001118),  
600 doi:[https://doi.org/10.1016/S0925-2312\(98\)00111-8](https://doi.org/10.1016/S0925-2312(98)00111-8).
- 601 Mei, S., Montanari, A., Nguyen, P.M., 2018. A mean field view  
602 of the landscape of two-layer neural networks. *Proceedings of the*  
603 *National Academy of Sciences* 115, E7665–E7671. URL: [https://](https://www.pnas.org/content/115/33/E7665)  
604 [www.pnas.org/content/115/33/E7665](https://www.pnas.org/content/115/33/E7665), doi:10.1073/pnas.1806579115,  
605 arXiv:<https://www.pnas.org/content/115/33/E7665.full.pdf>.
- 606 Montanelli, H., Du, Q., 2019. New error bounds for deep ReLU net-  
607 works using sparse grids. *SIAM Journal on Mathematics of Data Science*  
608 1, 78–92. URL: <https://doi.org/10.1137/18M1189336>, doi:10.1137/  
609 [18M1189336](https://doi.org/10.1137/18M1189336), arXiv:<https://doi.org/10.1137/18M1189336>.
- 610 Montanelli, H., Yang, H., 2020. Error bounds for deep ReLU net-  
611 works using the Kolmogorov-Arnold superposition theorem. *Neu-*  
612 *ral Networks* 129, 1 – 6. URL: [http://www.sciencedirect.com/](http://www.sciencedirect.com/science/article/pii/S0893608019304058)  
613 [science/article/pii/S0893608019304058](http://www.sciencedirect.com/science/article/pii/S0893608019304058), doi:[https://doi.org/10.](https://doi.org/10.1016/j.neunet.2019.12.013)  
614 [1016/j.neunet.2019.12.013](https://doi.org/10.1016/j.neunet.2019.12.013).
- 615 Montanelli, H., Yang, H., Du, Q., 2020. Deep ReLU networks overcome the  
616 curse of dimensionality for bandlimited functions. *Journal of Computa-*  
617 *tional Mathematics* .
- 618 Opschoor, J.A., Schwab, C., Zech, J., 2019. Exponential ReLU  
619 DNN expression of holomorphic maps in high dimension.  
620 <https://math.ethz.ch/sam/research/reports.html?id=839> 2019-35.
- 621 Petersen, P., Voigtlaender, F., 2018. Optimal approximation of piece-  
622 wise smooth functions using deep ReLU neural networks. *Neural*  
623 *Networks* 108, 296 – 330. URL: [http://www.sciencedirect.com/](http://www.sciencedirect.com/science/article/pii/S0893608018302454)  
624 [science/article/pii/S0893608018302454](http://www.sciencedirect.com/science/article/pii/S0893608018302454), doi:[https://doi.org/10.](https://doi.org/10.1016/j.neunet.2018.08.019)  
625 [1016/j.neunet.2018.08.019](https://doi.org/10.1016/j.neunet.2018.08.019).
- 626 Poggio, T., Mhaskar, H.N., Rosasco, L., Miranda, B., Liao, Q., 2017. Why  
627 and when can deep—but not shallow—networks avoid the curse of dimen-  
628 sionality: A review. *International Journal of Automation and Computing*  
629 14, 503–519.

630 Schmidt-Hieber, J., 2017. Nonparametric regression using deep neural net-  
631 works with ReLU activation function URL: [https://arxiv.org/abs/](https://arxiv.org/abs/1708.06633)  
632 [1708.06633](https://arxiv.org/abs/1708.06633).

633 Schmidt-Hieber, J., 2020. The kolmogorov-arnold representation theorem  
634 revisited. [arXiv:2007.15884](https://arxiv.org/abs/2007.15884).

635 Shen, Z., Yang, H., Zhang, S., 2019. Nonlinear approximation via composi-  
636 tions. Neural Networks 119, 74 – 84. URL: [http://www.sciencedirect.](http://www.sciencedirect.com/science/article/pii/S0893608019301996)  
637 [com/science/article/pii/S0893608019301996](http://www.sciencedirect.com/science/article/pii/S0893608019301996), doi:[https://doi.org/](https://doi.org/10.1016/j.neunet.2019.07.011)  
638 [10.1016/j.neunet.2019.07.011](https://doi.org/10.1016/j.neunet.2019.07.011).

639 Shen, Z., Yang, H., Zhang, S., 2020. Deep network approximation charac-  
640 terized by number of neurons. Communications in Computational Physics  
641 .

642 Shen, Z., Yang, H., Zhang, S., 2020. Deep network with approxima-  
643 tion error being reciprocal of width to power of square root of depth.  
644 arXiv:2006.12231 [arXiv:2006.12231](https://arxiv.org/abs/2006.12231).

645 Wu, L., Ma, C., E, W., 2018. How sgd selects the global minima in over-  
646 parameterized learning: A dynamical stability perspective, in: Bengio, S.,  
647 Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R.  
648 (Eds.), Advances in Neural Information Processing Systems 31. Curran  
649 Associates, Inc., pp. 8279–8288.

650 Yang, Y., Wang, Y., 2020. Approximation in shift-invariant  
651 spaces with deep ReLU neural networks. arXiv e-prints ,  
652 arXiv:2005.11949 [arXiv:2005.11949](https://arxiv.org/abs/2005.11949).

653 Yarotsky, D., 2017. Error bounds for approximations with deep ReLU  
654 networks. Neural Networks 94, 103 – 114. URL: [http://www.](http://www.sciencedirect.com/science/article/pii/S0893608017301545)  
655 [sciencedirect.com/science/article/pii/S0893608017301545](http://www.sciencedirect.com/science/article/pii/S0893608017301545),  
656 doi:<https://doi.org/10.1016/j.neunet.2017.07.002>.

657 Yarotsky, D., 2018. Optimal approximation of continuous functions by very  
658 deep ReLU networks, in: Bubeck, S., Perchet, V., Rigollet, P. (Eds.),  
659 Proceedings of the 31st Conference On Learning Theory, PMLR. pp. 639–  
660 649. URL: <http://proceedings.mlr.press/v75/yarotsky18a.html>.



661 Yarotsky, D., Zhevnerchuk, A., 2019. The phase diagram of ap-  
662 proximation rates for deep neural networks. arXiv e-prints ,  
663 arXiv:1906.09477[arXiv:1906.09477](https://arxiv.org/abs/1906.09477).