

OPTIMAL APPROXIMATION WITH SPARSELY CONNECTED DEEP NEURAL NETWORKS

HELMUT BÖLCSKEI*, PHILIPP GROHS†, GITTA KUTYNIOK‡, AND PHILIPP PETERSEN§

Abstract. We derive fundamental lower bounds on the connectivity and the memory requirements of deep neural networks guaranteeing uniform approximation rates for arbitrary function classes in $L^2(\mathbb{R}^d)$. In other words, we establish a connection between the complexity of a function class and the complexity of deep neural networks approximating functions from this class to within a prescribed accuracy. Additionally, we prove that our lower bounds are achievable for a broad family of function classes. Specifically, all function classes that are optimally approximated by a general class of representation systems—so-called *affine systems*—can be approximated by deep neural networks with minimal connectivity and memory requirements. Affine systems encompass a wealth of representation systems from applied harmonic analysis such as wavelets, ridgelets, curvelets, shearlets, α -shearlets, and more generally α -molecules. Our central result elucidates a remarkable universality property of neural networks and shows that they achieve the optimum approximation properties of all affine systems combined. As a specific example, we consider the class of α^{-1} -cartoon-like functions, which is approximated optimally by α -shearlets. We also explain how our results can be extended to the approximation of functions on low-dimensional immersed manifolds. Finally, we present numerical experiments demonstrating that the standard stochastic gradient descent algorithm yields deep neural networks with close-to-optimal approximation rates. Moreover, these results indicate that stochastic gradient descent can learn approximations that are sparse in the representation systems optimally sparsifying the function class the network is trained on.

Keywords. Neural networks, function approximation, optimal sparse approximation, sparse connectivity, wavelets, shearlets

AMS subject classification. 41A25, 82C32, 42C40, 42C15, 41A46, 68T05, 94A34, 94A12

1. Introduction. Neural networks arose from the seminal work by McCulloch and Pitts [41] in 1943 which, inspired by the functionality of the human brain, introduced an algorithmic approach to learning with the aim of building a theory of artificial intelligence. Roughly speaking, a neural network consists of neurons arranged in layers and connected by weighted edges; in mathematical terms this boils down to a concatenation of affine linear functions and relatively simple non-linearities.

Despite significant theoretical progress in the 1990s [12, 34], the area has seen practical progress only during the past decade, triggered by the drastic improvements in computing power and the availability of vast amounts of training data. Deep neural networks, i.e., networks with large numbers of layers, are now state-of-the-art technology for a wide variety of applications, such as image classification [36], speech recognition [33], or game intelligence [13]. For an in-depth overview, we refer to the survey paper by LeCun, Bengio, and Hinton [39] and the recent book [24].

A neural network effectively implements a non-linear mapping and can be used to either perform classification directly or to extract features that are then fed into a classifier, such as a support vector machine [54]. In the former case, the primary goal is to approximate an unknown classification function based on a given set of input-output value pairs. This is typically accomplished by learning the network’s weights through, e.g., the stochastic gradient descent (via backpropagation) algorithm [52]. In a classification task with, say, two classes,

*Department of Information Technology and Electrical Engineering and Department of Mathematics, ETH Zürich, 8092 Zürich, Switzerland. Email-Address: hboelcskei@ethz.ch

†Faculty of Mathematics, University of Vienna, 1090 Vienna, Austria, and Research Platform DataScience@UniVienna, University of Vienna, 1090 Vienna, Austria. Email-Address: philipp.grohs@univie.ac.at

‡Institut für Mathematik and Department of Electrical Engineering & Computer Science, Technische Universität Berlin, 10623 Berlin, Germany. Email-Address: kutyniok@math.tu-berlin.de

§Mathematical Institute, University of Oxford, OX26GG Oxford, United Kingdom. Email-Address: philipp.petersen@maths.ox.ac.uk

42 the function to be learned would take only two values, whereas in the case of, e.g., the pre-
 43 diction of the temperature in a certain environment, it would be real-valued. It is therefore
 44 clear that characterizing to what extent (deep) neural networks are capable of approximating
 45 general functions is a question of significant practical relevance.

46 Neural networks employed in practice often consist of hundreds of layers and may de-
 47 pend on billions of parameters, see for example the work [32] on image classification. Train-
 48 ing and operation of networks of this scale entail formidable computational challenges. As
 49 a case in point, we mention speech recognition on a smartphone such as, e.g., Apple’s SIRI-
 50 system, which operates in the cloud. Android’s speech recognition system has meanwhile
 51 released an offline version based on a neural network with sparse connectivity. The desire to
 52 reduce the complexity of network training and operation naturally leads to the question of the
 53 fundamental limits on function approximation through neural networks with sparse connec-
 54 tivity. In addition, the network’s memory requirements in terms of the number of bits needed
 55 to store its topology and weights are of concern in practice.

56 The purpose of this paper is to understand the connectivity and memory requirements
 57 of (deep) neural networks induced by demands on their approximation-theoretic properties.
 58 Specifically, defining the complexity of a function class \mathcal{C} as the rate of growth of the min-
 59 imum number of bits needed to describe any element in \mathcal{C} to within a maximum allowed
 60 error approaching zero, we shall be interested in the following question: Depending on the
 61 complexity of \mathcal{C} , what are the connectivity and memory requirements of a deep neural net-
 62 work approximating every element in \mathcal{C} to within an error of ε ? We address this question by
 63 interpreting the network as an encoder in Donoho’s min-max rate distortion theory [19] and
 64 establishing rate-distortion optimality for a broad family of function classes \mathcal{C} , namely those
 65 classes for which so-called affine systems—a general class of representation systems—yield
 66 optimal approximation rates in the sense of non-linear approximation theory [16]. Affine
 67 systems encompass a wealth of representation systems from applied harmonic analysis such
 68 as wavelets [14], ridgelets [4], curvelets [6], shearlets [31], α -shearlets and more generally
 69 α -molecules [27]. Our result therefore uncovers an interesting universality property of deep
 70 neural networks; they exhibit the optimal approximation properties of all affine systems com-
 71 bined. The technique we develop to prove our main statements is interesting in its own right
 72 as it constitutes a more general framework for transferring results on function approximation
 73 through representation systems to results on approximation by deep neural networks.

74 **1.1. Deep Neural Networks.** While various network architectures exist in the literature,
 75 we focus on the following setup.

76 **DEFINITION 1.1.** *Let $L, d, N_1, \dots, N_L \in \mathbb{N}$ with $L \geq 2$. A map $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^{N_L}$ given*
 77 *by*

$$78 \quad (1.1) \quad \Phi(x) = W_L \rho(W_{L-1} \rho(\dots \rho(W_1(x))))), \quad \text{for } x \in \mathbb{R}^d,$$

79 *with affine linear maps $W_\ell : \mathbb{R}^{N_{\ell-1}} \rightarrow \mathbb{R}^{N_\ell}$, $1 \leq \ell \leq L$, and the non-linear activation*
 80 *function ρ acting component-wise, is called a neural network. Here, $N_0 := d$ is the dimension*
 81 *of the 0-th layer referred to as the input layer, L denotes the number of layers (not counting*
 82 *the input layer), N_1, \dots, N_{L-1} stands for the dimensions of the $L - 1$ hidden layers, and*
 83 *N_L is the dimension of the output layer. The affine linear map W_ℓ is defined via $W_\ell(x) =$*
 84 *$A_\ell x + b_\ell$ with $A_\ell \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$ and the affine part $b_\ell \in \mathbb{R}^{N_\ell}$. $(A_\ell)_{i,j}$ represents the weight*
 85 *associated with the edge between the j -th node in the $(\ell - 1)$ -th layer and the i -th node in*
 86 *the ℓ -th layer, while $(b_\ell)_i$ is the weight associated with the i -th node in the ℓ -th layer. These*
 87 *assignments are schematized in Figure 1. The total number of nodes is given by $\mathcal{N}(\Phi) :=$*
 88 *$d + \sum_{\ell=1}^L N_\ell$. The real numbers $(A_\ell)_{i,j}$ and $(b_\ell)_i$ are said to be the network’s edge weights*

89 and node weights, respectively, and the total number of nonzero edge weights, denoted by
 90 $\mathcal{M}(\Phi)$, is the network's connectivity.

91 The term “network” stems from the interpretation of the mapping Φ as a weighted acyclic
 92 directed graph with nodes arranged in $L + 1$ hierarchical layers and edges only between ad-
 93 jacent layers. If the network's connectivity $\mathcal{M}(\Phi)$ is small relative to the number of connec-
 94 tions possible (i.e., the number of edges in the graph that is fully connected between adjacent
 95 layers), we say that the network is *sparsely connected*.

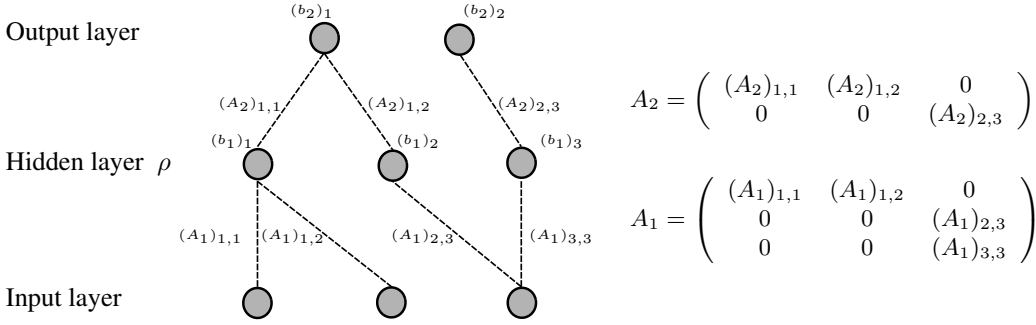


FIG. 1. Assignment of the weights $(A_\ell)_{i,j}$ and $(b_\ell)_i$ of a two-layer network to the edges and nodes, respectively.

96 Throughout the paper, we consider the case $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}$, i.e., $N_L = 1$, which includes
 97 situations such as the classification and temperature prediction problem described above. We
 98 emphasize, however, that the general results of Sections 3, 4, and 5 are readily generalized to
 99 $N_L > 1$.

100 We denote the class of networks $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}$ with exactly L layers, connectivity
 101 no more than M , and activation function ρ by $\mathcal{NN}_{L,M,d,\rho}$ with the understanding that for
 102 $L = 1$, the set $\mathcal{NN}_{L,M,d,\rho}$ is empty. Moreover, we let $\mathcal{NN}_{\infty,M,d,\rho} := \bigcup_{L \in \mathbb{N}} \mathcal{NN}_{L,M,d,\rho}$,
 103 $\mathcal{NN}_{L,\infty,d,\rho} := \bigcup_{M \in \mathbb{N}} \mathcal{NN}_{L,M,d,\rho}$, and $\mathcal{NN}_{\infty,\infty,d,\rho} := \bigcup_{L \in \mathbb{N}} \mathcal{NN}_{L,\infty,d,\rho}$.

104 Now, given a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, we are interested in the theoretically best possible
 105 approximation of f by a network $\Phi \in \mathcal{NN}_{\infty,M,d,\rho}$. Specifically, we will want to know how
 106 the approximation quality depends on the connectivity M and what the associated number
 107 of bits needed to store the network topology and the corresponding quantized weights is.
 108 Clearly, smaller M entails lower computational complexity in terms of evaluating (1.1) and
 109 a smaller number of bits translates to reduced memory requirements for storing the network.
 110 Such a result benchmarks all conceivable algorithms for learning the network topology and
 111 weights.

112 **1.2. Quantifying Approximation Quality.** We proceed to formalizing our problem
 113 statement and start with a brief review of a widely used framework in approximation the-
 114 ory [17, 16].

115 Fix $\Omega \subset \mathbb{R}^d$. Let \mathcal{C} be a compact set of functions in $L^2(\Omega)$, henceforth referred to as
 116 function class, and consider a corresponding system $\mathcal{D} := (\varphi_i)_{i \in I} \subset L^2(\Omega)$ with I count-
 117 able, termed *representation system*. We study the *error of best M -term approximation* of
 118 $f \in \mathcal{C}$ in \mathcal{D} :

119 **DEFINITION 1.2** ([17]). *Given $d \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$, a function class $\mathcal{C} \subset L^2(\Omega)$, and a*

120 representation system $\mathcal{D} = (\varphi_i)_{i \in I} \subset L^2(\Omega)$, we define, for $f \in \mathcal{C}$ and $M \in \mathbb{N}$,

$$121 \quad (1.2) \quad \Gamma_M^{\mathcal{D}}(f) := \inf_{\substack{I_M \subseteq I, \\ \#I_M = M, (c_i)_{i \in I_M}}} \left\| f - \sum_{i \in I_M} c_i \varphi_i \right\|_{L^2(\Omega)}.$$

123 We call $\Gamma_M^{\mathcal{D}}(f)$ the best M -term approximation error of f in \mathcal{D} . Every $f_M = \sum_{i \in I_M} c_i \varphi_i$
 124 attaining the infimum in (1.2) is referred to as a best M -term approximation of f in \mathcal{D} . The
 125 supremal $\gamma > 0$ such that

$$126 \quad \sup_{f \in \mathcal{C}} \Gamma_M^{\mathcal{D}}(f) \in \mathcal{O}(M^{-\gamma}), \quad M \rightarrow \infty,$$

127 will be denoted by $\gamma^*(\mathcal{C}, \mathcal{D})$. Here, $\mathcal{O}(g(\cdot))$ denotes the class of functions bounded asymptotically
 128 by g in the sense of standard Landau notation. We say that the best M -term approxi-
 129 mation rate of \mathcal{C} in the representation system \mathcal{D} is $\gamma^*(\mathcal{C}, \mathcal{D})$.

130 Function classes \mathcal{C} widely studied in the approximation theory literature include unit
 131 balls in Lebesgue, Sobolev, or Besov spaces [16], as well as α -cartoon-like functions [27].
 132 A wealth of structured representation systems \mathcal{D} is provided by the area of applied harmonic
 133 analysis, starting with wavelets [14], followed by ridgelets [4], curvelets [6], shearlets [31],
 134 parabolic molecules [29], and most generally α -molecules [27], which include all previously
 135 named systems as special cases. Further examples are Gabor frames [25] and wave atoms
 136 [15].

137 **1.3. Approximation by Deep Neural Networks.** The main conceptual contribution of
 138 this paper is the development of an approximation-theoretic framework for deep neural net-
 139 works in the spirit of [17]. Specifically, we shall substitute the concept of best M -term
 140 approximation with representation systems by best M -edge approximation through neural
 141 networks. In other words, parsimony in terms of the number of participating elements of a
 142 representation system is replaced by parsimony in terms of connectivity. More formally, we
 143 consider the following setup.

144 **DEFINITION 1.3.** Given $d \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$, a function class $\mathcal{C} \subset L^2(\Omega)$, and an activation
 145 function $\rho : \mathbb{R} \rightarrow \mathbb{R}$, we define, for $f \in \mathcal{C}$ and $M \in \mathbb{N}$,

$$146 \quad (1.3) \quad \Gamma_M^{\mathcal{NN}}(f) := \inf_{\Phi \in \mathcal{NN}_{\infty, M, d, \rho}} \|f - \Phi\|_{L^2(\Omega)}.$$

148 We call $\Gamma_M^{\mathcal{NN}}(f)$ the best M -edge approximation error of f . The supremal $\gamma > 0$ such that

$$149 \quad \sup_{f \in \mathcal{C}} \Gamma_M^{\mathcal{NN}}(f) \in \mathcal{O}(M^{-\gamma}), \quad M \rightarrow \infty,$$

150 will be denoted by $\gamma_{\mathcal{NN}}^*(\mathcal{C}, \rho)$. We say that the best M -edge approximation rate of \mathcal{C} by neural
 151 networks with activation function ρ is $\gamma_{\mathcal{NN}}^*(\mathcal{C}, \rho)$.

152 We emphasize that the infimum in (1.3) is taken over all networks with fixed activation
 153 function ρ , fixed input dimension d , no more than M edges of nonzero weight, and arbitrary
 154 number of layers L . In particular, this means that the infimum is taken over all possible net-
 155 work topologies. The resulting best M -edge approximation rate is fundamental as it bench-
 156 marks all learning algorithms, i.e., all algorithms that map an input function f and an $\varepsilon > 0$
 157 to a neural network that approximates f with error no more than ε . Our framework hence
 158 provides a means for assessing the performance of a given learning algorithm in the sense of
 159 allowing to measure how close the M -edge approximation rate induced by the algorithm is
 160 to the best M -edge approximation rate $\gamma_{\mathcal{NN}}^*(\mathcal{C}, \rho)$.

161 **1.4. Previous Work.** The best-known results on approximation by neural networks are
 162 the universal approximation theorems of Hornik [34] and Cybenko [12], stating that every
 163 measurable function f can be approximated arbitrarily well by a single-hidden-layer ($L = 2$
 164 in our terminology) neural network. The literature on approximation-theoretic properties
 165 of networks with a single hidden layer continuing this line of work is abundant. Without
 166 any claim to completeness, we mention work on approximation error bounds in terms of the
 167 number of neurons for functions with bounded first moments [1], [2], the non-existence of
 168 localized approximations [7], a fundamental lower bound on approximation rates [18, 4], and
 169 the approximation of smooth or analytic functions [44, 42].

170 Approximation-theoretic results for networks with multiple hidden layers were obtained
 171 in [35, 43] for general functions, in [23] for continuous functions, and for functions to-
 172 gether with their derivatives in [47]. In [7] it was shown that for certain approximation
 173 tasks deep networks can perform fundamentally better than single-hidden-layer networks.
 174 We also highlight two recent papers, which investigate the benefit—from an approximation-
 175 theoretic perspective—of multiple hidden layers. Specifically, in [21] it was shown that there
 176 exists a function which, although expressible through a small three-layer network, can only
 177 be represented through a very large two-layer network; here size is measured in terms of
 178 the total number of neurons in the network. In the setting of deep neural networks, first re-
 179 sults of a nature similar to those in [21] were reported in [46]. For the activation function
 180 $\rho(x) = \max\{0, x\}$ —usually referred to as ReLU—it was demonstrated in [58, 50] that deep
 181 networks approximate sufficiently smooth functions with rates higher than those achieved by
 182 shallow networks.

183 Convolutional neural networks are obtained as a special case of the general networks
 184 considered in this paper, namely by taking the affine linear transformations in Definition 1.1
 185 to have circulant A -matrices. Linking the expressivity properties of neural networks to tensor
 186 decompositions, [9, 10] establish the existence of functions that can be realized by relatively
 187 small deep convolutional networks but require exponentially larger shallow convolutional net-
 188 works. Universal approximation theorems for convolutional neural networks are provided in
 189 [59, 57]. The approximation-theoretic equivalence results between convolutional networks
 190 and general networks established in [49], together with the main findings of the present pa-
 191 per, lead to upper and lower bounds on approximation rates attainable with convolutional
 192 networks.

193 For survey articles on approximation-theoretic aspects of neural networks, we refer the
 194 interested reader to [22, 51].

195 Most closely related to our work is that by Shaham, Cloninger, and Coifman [53], which
 196 shows that for functions that are sparse in specific wavelet frames, the best M -edge approxi-
 197 mation rate of three-layer neural networks is at least as high as the best M -term approximation
 198 rate in piecewise linear wavelet frames.

199 **1.5. Contributions.** Our contributions can be grouped into four threads.

- 200 • *Fundamental lower bound on connectivity.* We quantify the minimum network con-
 201 nectivity needed to allow approximation of *all* elements of a given function class \mathcal{C}
 202 to within a maximum allowed error. On a conceptual level, this result establishes a
 203 universal link between the complexity of a given function class and the connectivity
 204 required by corresponding approximating neural networks.
- 205 • *Transfer from M -term to M -edge approximation.* We develop a general framework
 206 for transferring best M -term approximation results in representation systems to best
 207 M -edge approximation results for neural networks.

- 208 • *Memory requirements.* We characterize the memory requirements needed to store
209 the topology and the quantized weights of optimally-approximating neural networks.
- 210 • *Realizability of optimal approximation rates.* An important practical question is how
211 neural networks trained by stochastic gradient descent (via backpropagation) [52]
212 perform relative to the fundamental bounds established in the paper. Interestingly,
213 our numerical experiments indicate that stochastic gradient descent can achieve M -
214 edge approximation rates quite close to the fundamental limit.

215 **1.6. Outline of the Paper.** Section 2 introduces the novel concept of effective best M -
216 edge approximation. The fundamental lower bound on connectivity is developed in Section 3.
217 Section 4 describes a general framework for transferring best M -term approximation results
218 in representation systems to best M -edge approximation results for neural networks. In Section
219 5, we apply this transfer framework to the broad class of affine representation systems,
220 and Section 6 shows that this leads to optimal M -edge approximation rates for cartoon functions.
221 In Section 7, we briefly outline the extension of our main findings to the approximation
222 of functions defined on manifolds. Finally, numerical results assessing the performance of
223 stochastic gradient descent (via backpropagation) relative to our lower bound on connectivity
224 are reported in Section 8.

225 **2. Effective Best M -term and Best M -edge Approximation.** We proceed by introduc-
226 ing M -term approximation via dictionaries and M -edge approximation via neural networks.
227 These concepts do, however, not allow for a meaningful notion of optimality in practice. A
228 remedy is provided by effective best M -term approximation according to [19, 26] and the
229 new concept of effective best M -edge approximation introduced below.

230 **2.1. Effective Best M -term Approximation.** The best M -term approximation rate
231 $\gamma^*(\mathcal{C}, \mathcal{D})$ according to Definition 1.2 measures the hardness of approximation of a given
232 function class \mathcal{C} by a fixed representation system \mathcal{D} . It is sensible to ask whether for a given
233 function class \mathcal{C} , there is a fundamental limit on $\gamma^*(\mathcal{C}, \mathcal{D})$ when one is allowed to vary over
234 \mathcal{D} . As shown in [19, 26], every dense (and countable) $\mathcal{D} \subset L^2(\Omega)$, $\Omega \subset \mathbb{R}^d$, results in
235 $\gamma^*(\mathcal{C}, \mathcal{D}) = \infty$ for all function classes $\mathcal{C} \subset L^2(\Omega)$. However, identifying the elements in \mathcal{D}
236 participating in the best M -term approximation is infeasible as it entails searching through
237 the infinite set \mathcal{D} and requires, in general, an infinite number of bits to describe the indices of
238 the participating elements. This insight leads to the concept of “best M -term approximation
239 subject to polynomial-depth search” as introduced by Donoho in [19].

240 **DEFINITION 2.1.** Given $d \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$, a function class $\mathcal{C} \subset L^2(\Omega)$, and a represen-
241 tation system $\mathcal{D} = (\varphi_i)_{i \in I} \subset L^2(\Omega)$, the supremal $\gamma > 0$ so that there exist a polynomial π
242 and a constant $D > 0$ such that

$$243 \quad (2.1) \quad \sup_{f \in \mathcal{C}} \inf_{\substack{I_M \subset \{1, \dots, \pi(M)\}, \\ \#I_M = M, (c_i)_{i \in I_M}, \max_{i \in I_M} |c_i| \leq D}} \left\| f - \sum_{i \in I_M} c_i \varphi_i \right\|_{L^2(\Omega)} \in \mathcal{O}(M^{-\gamma}), \quad M \rightarrow \infty,$$

245 will be denoted by $\gamma^{*,\text{eff}}(\mathcal{C}, \mathcal{D})$ and referred to as effective best M -term approximation rate
246 of \mathcal{C} in the representation system \mathcal{D} .

247 We will demonstrate in Section 3.2 that $\sup_{\mathcal{D} \subset L^2(\Omega)} \gamma^{*,\text{eff}}(\mathcal{C}, \mathcal{D})$ is, indeed, finite under
248 quite general conditions on \mathcal{C} and, in particular, depends on the “description complexity” of
249 \mathcal{C} . This will allow us to assess the approximation capacity of a given representation system
250 \mathcal{D} for \mathcal{C} by comparing $\gamma^{*,\text{eff}}(\mathcal{C}, \mathcal{D})$ to the ultimate limit $\sup_{\mathcal{D} \subset L^2(\Omega)} \gamma^{*,\text{eff}}(\mathcal{C}, \mathcal{D})$.

251 **2.2. Effective Best M -edge Approximation.** We next aim at establishing a relation-
252 ship in the spirit of effective best M -term approximation for approximation through deep

253 neural networks. To this end, we first note that Definition 1.3 encounters problems sim-
 254 ilar to those identified for approximation by representation systems, namely the quantity
 255 $\sup_{\rho: \mathbb{R} \rightarrow \mathbb{R}} \gamma_{\mathcal{NN}}^*(\mathcal{C}, \rho)$ does not reveal anything tangible about the approximation complex-
 256 ity of \mathcal{C} in deep neural networks, unless further constraints are imposed on the approximating
 257 network. To make this point, we first review the following remarkable result:

258 **THEOREM 2.2** ([40]). *There exists a function $\rho : \mathbb{R} \rightarrow \mathbb{R}$ that is C^∞ , strictly increasing,
 259 and satisfies $\lim_{x \rightarrow \infty} \rho(x) = 1$ and $\lim_{x \rightarrow -\infty} \rho(x) = 0$, such that for any $d \in \mathbb{N}$, any
 260 $f \in C([0, 1]^d)$, and any $\varepsilon > 0$, there is a neural network Φ with activation function ρ and
 261 three layers, of dimensions $N_1 = 3d$, $N_2 = 6d + 3$, and $N_3 = 1$, satisfying*

$$262 \quad (2.2) \quad \sup_{x \in [0, 1]^d} |f(x) - \Phi(x)| \leq \varepsilon.$$

263 We observe that the number of nodes and the number of layers of the approximating
 264 network in Theorem 2.2 do not depend on the approximation error ε . In particular, ε can be
 265 chosen arbitrarily small while having $\mathcal{M}(\Phi)$ bounded. By density of $C([0, 1]^d)$ in $L^2([0, 1]^d)$
 266 and hence in all compact subsets of $L^2([0, 1]^d)$, this implies the existence of an activation
 267 function $\rho : \mathbb{R} \rightarrow \mathbb{R}$ such that $\gamma_{\mathcal{NN}}^*(\mathcal{C}, \rho) = \infty$ for all compact $\mathcal{C} \subset L^2([0, 1]^d)$, $d \in$
 268 \mathbb{N} . However, the networks underlying Theorem 2.2 necessarily lead to weights that are (in
 269 absolute value) not bounded by $|\pi(\varepsilon^{-1})|$ for a polynomial π , a requirement we will have to
 270 impose to get rate-distortion-optimal approximation through neural networks (see Section 3).
 271 To see that the weights, indeed, do not obey a polynomial growth bound in ε^{-1} , we note that
 272 thanks to Theorem 2.2, there exist $C > 0$ and $\gamma > 0$ such that

$$273 \quad (2.3) \quad \sup_{f \in \mathcal{C}} \inf_{\Phi_M \in \mathcal{NN}_{3, M, d, \rho}} \|f - \Phi_M\|_{L^2(\Omega)} \leq CM^{-\gamma}, \text{ for all } M \in \mathbb{N}.$$

275 Now, as ε in Theorem 2.2 can be made arbitrarily small while the connectivity of the corre-
 276 sponding networks remains upper-bounded by $21d^2 + 15d + 3$, (2.3) would have to hold for
 277 arbitrarily large γ , in particular also for $\gamma > \gamma_{\mathcal{NN}}^{*, \text{eff}}(\mathcal{C}, \rho)$, where $\gamma_{\mathcal{NN}}^{*, \text{eff}}(\mathcal{C}, \rho)$ is the effective
 278 best M -edge approximation rate according to Definition 2.3. By Theorem 3.4 below, how-
 279 ever, $\gamma_{\mathcal{NN}}^{*, \text{eff}}(\mathcal{C}, \rho) \leq \gamma^*(\mathcal{C})$, where $\gamma^*(\mathcal{C})$ is the optimal exponent according to Definition 3.1.
 280 Owing to Definition 2.3, we can therefore conclude that the weights of the network achieving
 281 the infimum in (2.3) can not be bounded by a polynomial in $M \sim \varepsilon^{-1}$ whenever $\gamma^*(\mathcal{C}) < \infty$.
 282 Here and in the sequel, we write $a \sim b$ if the variables a and b are proportional, i.e., there
 283 exist uniform constants $c_1, c_2 > 0$ such that $c_1 a \leq b \leq c_2 a$.

284 The observation just made resembles the problem in best M -term approximation which
 285 eventually led to the concept of *effective* best M -term approximation, where we restricted
 286 the search depth in the representation system \mathcal{D} to be polynomially bounded in M and the
 287 coefficients c_i to be bounded according to $\max_{i \in I_M} |c_i| \leq D$. Interpreting the weights
 288 in the network as the counterpart of the coefficients c_i in best M -term approximation, we
 289 see that the restriction on the search depth corresponds to restricting the size of the indices
 290 enumerating the participating weights. The need for such a restriction is obviated by the tree
 291 structure of deep neural networks as exposed in detail in the proof of Proposition 3.6. The
 292 second restriction will lead us to a growth condition on the weights, which is more generous
 293 than the corresponding requirement of the c_i in effective best M -term approximation being
 294 bounded.

295 In summary, this leads to the novel concept of “best M -edge approximation subject to
 296 polynomial weight growth” as formalized next.

297 **DEFINITION 2.3.** *Given $d \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$, a function class $\mathcal{C} \subset L^2(\Omega)$, and an activation
 298 function $\rho : \mathbb{R} \rightarrow \mathbb{R}$, the supremal $\gamma > 0$ so that there exist an $L \in \mathbb{N}$ and a polynomial π*

299 *such that*

$$300 \quad (2.4) \quad \sup_{f \in \mathcal{C}} \inf_{\Phi_M \in \mathcal{NN}_{L,M,d,\rho}^\pi} \|f - \Phi_M\|_{L^2(\Omega)} \in \mathcal{O}(M^{-\gamma}), \quad M \rightarrow \infty,$$

302 *where $\mathcal{NN}_{L,M,d,\rho}^\pi$ denotes the class of networks in $\mathcal{NN}_{L,M,d,\rho}$ that have all their weights*
 303 *bounded in absolute value by $|\pi(M)|$, will be referred to as effective best M -edge approxi-*
 304 *mation rate of \mathcal{C} by neural networks and denoted by $\gamma_{\mathcal{NN}}^{*,\text{eff}}(\mathcal{C}, \rho)$.*

305 We will show in Corollary 3.4 that $\sup_{\rho: \mathbb{R} \rightarrow \mathbb{R}} \gamma_{\mathcal{NN}}^{*,\text{eff}}(\mathcal{C}, \rho)$ is bounded and depends on the
 306 “description complexity” of the function class \mathcal{C} .

3. Fundamental Bounds on Effective M -Term and M -Edge Approximation Rate.

308 The purpose of this section is to establish fundamental bounds on effective best M -term and
 309 effective best M -edge approximation rates by evaluating the corresponding approximation
 310 strategies in the min-max rate distortion theory framework as developed in [19, 26].

311 **3.1. Min-Max Rate Distortion Theory.** Min-max rate distortion theory provides a the-
 312 oretical foundation for deterministic lossy data compression. We recall the following notions
 313 and concepts from [19, 26].

314 Let $d \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$, and consider the function class $\mathcal{C} \subset L^2(\Omega)$. Then, for each $\ell \in \mathbb{N}$,
 315 we denote by

$$316 \quad \mathfrak{E}^\ell := \{E : \mathcal{C} \rightarrow \{0, 1\}^\ell\}$$

318 the set of *binary encoders of \mathcal{C} of length ℓ* , and we let

$$319 \quad \mathfrak{D}^\ell := \{D : \{0, 1\}^\ell \rightarrow L^2(\Omega)\}$$

321 be the set of *binary decoders of length ℓ* . An encoder-decoder pair $(E, D) \in \mathfrak{E}^\ell \times \mathfrak{D}^\ell$ is said
 322 to *achieve uniform error ε over the function class \mathcal{C}* , if

$$323 \quad \sup_{f \in \mathcal{C}} \|D(E(f)) - f\|_{L^2(\Omega)} \leq \varepsilon.$$

325 A quantity of central interest is the minimal length $\ell \in \mathbb{N}$ for which there exists an
 326 encoder-decoder pair $(E, D) \in \mathfrak{E}^\ell \times \mathfrak{D}^\ell$ that achieves uniform error ε over the function class
 327 \mathcal{C} , along with its asymptotic behavior as made precise in the following definition.

328 **DEFINITION 3.1.** *Let $d \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$, and $\mathcal{C} \subset L^2(\Omega)$. Then, for $\varepsilon > 0$, the minimax*
 329 *code length $L(\varepsilon, \mathcal{C})$ is*

$$330 \quad L(\varepsilon, \mathcal{C}) := \min \left\{ \ell \in \mathbb{N} : \exists (E, D) \in \mathfrak{E}^\ell \times \mathfrak{D}^\ell : \sup_{f \in \mathcal{C}} \|D(E(f)) - f\|_{L^2(\Omega)} \leq \varepsilon \right\}.$$

331 *Moreover, the optimal exponent $\gamma^*(\mathcal{C})$ is defined as*

$$332 \quad \gamma^*(\mathcal{C}) := \sup \left\{ \gamma \in \mathbb{R} : L(\varepsilon, \mathcal{C}) \in \mathcal{O}\left(\varepsilon^{-\frac{1}{\gamma}}\right), \varepsilon \rightarrow 0 \right\}.$$

333 The optimal exponent $\gamma^*(\mathcal{C})$ quantifies the minimum growth rate of $L(\varepsilon, \mathcal{C})$ as the error ε
 334 tends to zero and can hence be seen as quantifying the “description complexity” of the func-
 335 tion class \mathcal{C} . Larger $\gamma^*(\mathcal{C})$ results in smaller growth rate and hence smaller memory require-
 336 ments for storing signals $f \in \mathcal{C}$ such that reconstruction with uniformly bounded error is
 337 possible. The quantity $\gamma^*(\mathcal{C})$ is closely related to the concept of Kolmogorov entropy [48].
 338 Remark 5.10 in [26] makes this connection explicit.

339 The optimal exponent is known for several function classes, such as subsets of Besov
 340 spaces $B_{p,q}^s(\mathbb{R}^d)$ with $1 \leq p, q < \infty, s > 0$, and $q > (s + 1/2)^{-1}$, namely all functions
 341 in $B_{p,q}^s(\mathbb{R}^d)$ of bounded norm, see e.g. [8]. If \mathcal{C} is a bounded subset of $B_{p,q}^s(\mathbb{R}^d)$, then
 342 we have $\gamma^*(\mathcal{C}) = s/d$. In the present paper, we shall be particularly interested in so-called
 343 β -cartoon-like functions, for which the optimal exponent is given by $\beta/2$ (see [20, 28] and
 344 Theorem 6.3).

345 **3.2. Fundamental Bound on Effective Best M -Term Approximation Rate.** We next
 346 recall a result from [19, 26], which says that, for a given function class \mathcal{C} , the optimal exponent
 347 $\gamma^*(\mathcal{C})$ constitutes a fundamental bound on the effective best M -term approximation rate of \mathcal{C}
 348 in any representation system. This gives operational meaning to $\gamma^*(\mathcal{C})$.

THEOREM 3.2 ([19, 26]). *Let $d \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$, $\mathcal{C} \subset L^2(\Omega)$, and assume that the effective best M -term approximation rate of \mathcal{C} in $\mathcal{D} \subset L^2(\Omega)$ is $\gamma^{*,\text{eff}}(\mathcal{C}, \mathcal{D})$. Then, we have*

$$\gamma^{*,\text{eff}}(\mathcal{C}, \mathcal{D}) \leq \gamma^*(\mathcal{C}).$$

349 In light of this result the following definition is natural (see also [26]).

DEFINITION 3.3. *Let $d \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$, and assume that the effective best M -term approximation rate of $\mathcal{C} \subset L^2(\Omega)$ in $\mathcal{D} \subset L^2(\Omega)$ is $\gamma^{*,\text{eff}}(\mathcal{C}, \mathcal{D})$. If*

$$\gamma^{*,\text{eff}}(\mathcal{C}, \mathcal{D}) = \gamma^*(\mathcal{C}),$$

350 *then the function class \mathcal{C} is said to be optimally representable by \mathcal{D} .*

351 **3.3. Fundamental Bound on Effective Best M -Edge Approximation Rate.** We now
 352 state the first main result of the paper, namely the equivalent of Theorem 3.2 for approxi-
 353 mation by deep neural networks. Specifically, we establish that the optimal exponent $\gamma^*(\mathcal{C})$
 354 also constitutes a fundamental bound on the effective best M -edge approximation rate of
 355 \mathcal{C} . We say below that a function $f : \mathbb{R} \rightarrow \mathbb{R}$ is *dominated* by a function $g : \mathbb{R} \rightarrow \mathbb{R}$ if
 356 $|f(x)| \leq |g(x)|$, for all $x \in \mathbb{R}$.

THEOREM 3.4. *Let $d \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$ be bounded, and $\mathcal{C} \subset L^2(\Omega)$. Then, for all $\rho : \mathbb{R} \rightarrow \mathbb{R}$ that are Lipschitz-continuous or differentiable with ρ' dominated by a polynomial, we have*

$$\gamma_{\mathcal{NN}}^{*,\text{eff}}(\mathcal{C}, \rho) \leq \gamma^*(\mathcal{C}).$$

357 The key ingredients of the proof of Theorem 3.4 are developed throughout this section
 358 and the formal proof will be stated at the end of the section. Before embarking, we note that,
 359 in analogy to Definition 3.3, what we just found suggests the following.

DEFINITION 3.5. *For $d \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$ bounded, we say that the function class $\mathcal{C} \subset L^2(\Omega)$ is optimally representable by neural networks with activation function $\rho : \mathbb{R} \rightarrow \mathbb{R}$, if*

$$\gamma_{\mathcal{NN}}^{*,\text{eff}}(\mathcal{C}, \rho) = \gamma^*(\mathcal{C}).$$

360 It is remarkable that the fundamental limits of approximation through representation sys-
 361 tems and approximation through deep neural networks are determined by the same quantity,
 362 although the approximants in the two cases are vastly different, namely linear combinations
 363 of elements of a representation system with the participating functions identified subject to a
 364 polynomial-depth search constraint in the former, and concatenations of affine functions fol-
 365 lowed by non-linearities under growth constraints on the weights in the network in the latter
 366 case.

367 A key ingredient of the proof of Theorem 3.4 is the following result, which establishes a
 368 fundamental lower bound on the connectivity of networks with quantized weights achieving
 369 uniform error ε over a given function class.

370 **PROPOSITION 3.6.** *Let $d \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$, $\rho : \mathbb{R} \rightarrow \mathbb{R}$, $c > 0$, and $\mathcal{C} \subset L^2(\Omega)$. Further, let*

$$371 \quad \mathbf{Learn} : \left(0, \frac{1}{2}\right) \times \mathcal{C} \rightarrow \mathcal{NN}_{\infty, \infty, d, \rho}$$

372 *be a map such that, for each pair $(\varepsilon, f) \in (0, 1/2) \times \mathcal{C}$, every weight of the neural network*
 373 *$\mathbf{Learn}(\varepsilon, f)$ is represented by no more than $\lceil c \log_2(\varepsilon^{-1}) \rceil$ bits* while guaranteeing that*

$$374 \quad (3.1) \quad \sup_{f \in \mathcal{C}} \|f - \mathbf{Learn}(\varepsilon, f)\|_{L^2(\Omega)} \leq \varepsilon.$$

375 *Then,*

$$376 \quad (3.2) \quad \sup_{f \in \mathcal{C}} \mathcal{M}(\mathbf{Learn}(\varepsilon, f)) \notin \mathcal{O}\left(\varepsilon^{-\frac{1}{\gamma}}\right), \varepsilon \rightarrow 0, \quad \text{for all } \gamma > \gamma^*(\mathcal{C}).$$

377 *Proof.* The proof is by contradiction. To this end, let $\gamma > \gamma^*(\mathcal{C})$ and assume that
 378 $\sup_{f \in \mathcal{C}} \mathcal{M}(\mathbf{Learn}(\varepsilon, f)) \in \mathcal{O}(\varepsilon^{-1/\gamma})$, $\varepsilon \rightarrow 0$. The contradiction will be effected by
 379 constructing encoder-decoder pairs $(E_\varepsilon, D_\varepsilon) \in \mathfrak{E}^{\ell(\varepsilon)} \times \mathfrak{D}^{\ell(\varepsilon)}$ achieving uniform error ε over
 380 \mathcal{C} with

$$381 \quad (3.3) \quad \ell(\varepsilon) \leq C_0 \cdot \sup_{f \in \mathcal{C}} [\mathcal{M}(\mathbf{Learn}(\varepsilon, f)) \log_2(\mathcal{M}(\mathbf{Learn}(\varepsilon, f))) + 1] \log_2(\varepsilon^{-1})$$

$$382 \quad \leq C_0 \cdot \left[\varepsilon^{-\frac{1}{\gamma}} \log_2\left(\varepsilon^{-\frac{1}{\gamma}}\right) + 1 \right] \log_2(\varepsilon^{-1})$$

$$383 \quad \leq C_1 \left(\varepsilon^{-\frac{1}{\gamma}} (\log_2(\varepsilon^{-1}))^2 + \log_2(\varepsilon^{-1}) \right) \in \mathcal{O}\left(\varepsilon^{-\frac{1}{\nu}}\right), \text{ for } \varepsilon \rightarrow 0,$$

385 where $C_0, C_1 > 0$ are constants and $\gamma > \nu > \gamma^*(\mathcal{C})$. Such a construction stands in contra-
 386 diction to the optimality of $\gamma^*(\mathcal{C})$ according to Definition 3.1.

387 We proceed to the construction of the encoder-decoder pairs $(E_\varepsilon, D_\varepsilon) \in \mathfrak{E}^{\ell(\varepsilon)} \times \mathfrak{D}^{\ell(\varepsilon)}$.
 388 Fix $f \in \mathcal{C}$. We enumerate the nodes in $\mathbf{Learn}(\varepsilon, f)$ by assigning natural numbers, henceforth
 389 called *indices*, increasing from left to right in every layer as schematized in Figure 2. For the
 390 sake of notational simplicity, we also set $\Phi := \mathbf{Learn}(\varepsilon, f)$ and $M := \mathcal{M}(\Phi)$. Without
 391 loss of generality, we henceforth assume that M is a power of 2 and larger than 1. The case
 392 $M = 0$ will be dealt with in Step 1 below. For all M that are not powers of 2 and for $M = 1$,
 393 we make use of the fact that $\mathcal{NN}_{L, M, d, \rho} \subset \mathcal{NN}_{L, M', d, \rho}$, where M' is the smallest power of 2
 394 larger than M , and we encode the network like an M' -edge network. Since $M < M' \leq 2M$
 395 this affects $\ell(\varepsilon)$ by a multiplicative constant only.

396 We recall that the number of layers of Φ is denoted by L , the number of nodes in these
 397 layers is N_1, \dots, N_L (see Definition 1.1), and d stands for the dimension of the input layer.

398 Denoting the number of nodes in layer $\ell = 1, \dots, L-1$ associated with edges of nonzero
 399 weight in the following layer by \tilde{N}_ℓ and setting $\tilde{N}_L = N_L = 1$, it follows that

$$400 \quad (3.4) \quad d + \sum_{\ell=1}^L \tilde{N}_\ell \leq 2\tilde{M},$$

402 where we let $\tilde{M} := M + d$. All other nodes do not contribute to the mapping $\Phi(x)$ and can
 403 hence be ignored. Moreover, we can assume that

$$404 \quad (3.5) \quad L \leq \tilde{M}$$

*Throughout the paper, we say that a weight of $\mathbf{Learn}(\varepsilon, f)$ is represented by no more than K bits, if it is taken from a set that is independent of f and has cardinality at most 2^K .

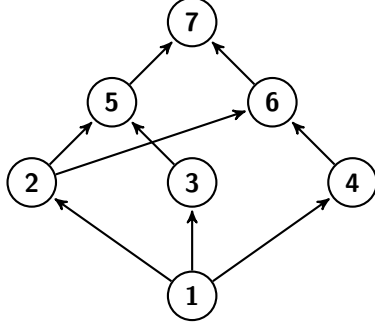


FIG. 2. Enumeration of nodes as employed in the proof of the theorem.

406 as otherwise there would be at least one layer $\ell > 1$ such that

~~407~~

$$A_\ell = 0.$$

409 As a consequence, the reduced network

~~410~~

$$x \mapsto W_L \rho(W_{L-1} \dots W_{\ell+1} \rho(0 \cdot x + b_\ell))$$

412 realizes the same function as the original network Φ but has fewer than L layers. This reduc-
413 tion can be repeated inductively until the resulting reduced network satisfies (3.5).

414 The bitstring representing Φ is constructed according to the following steps.

415 *Step 1:* If $M = 0$, we encode the network by a leading 0 followed by the bitstring
416 representing the node weight in the last layer. Upon defining $0 \log_2(0) = 0$, we then note that
417 (3.3) holds trivially and we terminate the encoding procedure. Else, we encode the number
418 of nonzero edge weights, M , by starting the overall bitstring with M 1's followed by a single
419 0. The length of this bitstring is therefore bounded by \widetilde{M} .

420 *Step 2:* We continue by encoding the number of layers in the network. Thanks to (3.5)
421 this requires no more than $\log_2(\widetilde{M})$ bits. We thus reserve the next $\log_2(\widetilde{M})$ bits for the binary
422 representation of L .

423 *Step 3:* Next, we store the dimension d of the input layer and the numbers of nodes
424 $\widetilde{N}_\ell, \ell = 1, \dots, L$, associated with edges of nonzero weight. As by (3.4) $d \leq \widetilde{M}$ and $\widetilde{N}_\ell \leq$
425 $2\widetilde{M}$, for all ℓ , we can encode (generously) d and each \widetilde{N}_ℓ using $\log_2(\widetilde{M}) + 1$ bits. For the
426 sake of concreteness, we first encode d followed by $\widetilde{N}_1, \dots, \widetilde{N}_L$ in that order. In total, Step 3
427 requires a bitstring of length

~~428~~
~~429~~

$$\left((L+1) \cdot \left(\log_2(\widetilde{M}) + 1 \right) \right) \leq (\widetilde{M} + 1) \log_2(\widetilde{M}) + \widetilde{M} + 1.$$

430 In combination with Steps 1 and 2 this yields an overall bitstring of length at most

~~431~~

$$(3.6) \quad \widetilde{M} \log_2(\widetilde{M}) + 2 \log_2(\widetilde{M}) + 2\widetilde{M} + 1.$$

432 *Step 4:* We encode the topology of the graph associated with Φ and consider only nodes
433 that contribute to the mapping $\Phi(x)$. Recall that we assigned a unique index i , ranging from 1
434 to $\widetilde{N} := d + \sum_{\ell=1}^L \widetilde{N}_\ell$, to each of these nodes. By (3.4) each of these indices can be encoded
435 by a bitstring of length $\log_2(\widetilde{M}) + 1$. We denote the bitstring corresponding to index i by
436 $b(i) \in \{0, 1\}^{\log_2(\widetilde{M})+1}$ and let $n(i)$ be the number of children of the node with index i , i.e.,

437 the number of nodes in the next layer connected to the node with index i via an edge. (Here,
 438 $n(\tilde{N}) = 0$.) For each node $i = 1, \dots, \tilde{N}$, we form a bitstring of length $n(i) \cdot (\log_2(\tilde{M}) + 1)$
 439 by concatenating the bitstrings $b(j)$ for all j such that there is an edge between i and j . We
 440 follow this string with an all-zeros bitstring of length $\log_2(\tilde{M}) + 1$ to signal the transition to
 441 the node with index $i + 1$. The enumeration is concluded with an all-zeros bitstring of length
 442 $\log_2(\tilde{M}) + 1$ signaling that the last node has been reached. Overall, this yields a bitstring of
 443 length

$$444 \quad (3.7) \quad \sum_{i=1}^{\tilde{N}} (n(i) + 1) \cdot (\log_2(\tilde{M}) + 1) \leq 3\tilde{M} \cdot (\log_2(\tilde{M}) + 1),$$

445 where we used $\sum_{i=1}^{\tilde{N}} n(i) = M < \tilde{M}$ and (3.4). Combining (3.6) and (3.7) it follows that
 446 we have encoded the overall topology of the network Φ using at most

$$447 \quad (3.8) \quad 5\tilde{M} + 4\tilde{M} \log_2(\tilde{M}) + 2 \log_2(\tilde{M}) + 1$$

448 bits.

449 *Step 5:* We encode the weights of Φ . By assumption, each weight can be represented
 450 by a bitstring of length $\lceil c \log_2(\varepsilon^{-1}) \rceil$. For each node $i = 1, \dots, \tilde{N}$, we reserve the first
 451 $\lceil c \log_2(\varepsilon^{-1}) \rceil$ bits to encode its associated node weight and, for each of its children a bitstring
 452 of length $\lceil c \log_2(\varepsilon^{-1}) \rceil$ to encode the weight corresponding to the edge between that child
 453 and its parent node. Concatenating the results in ascending order of child node indices, we
 454 get a bitstring of length $(n(i) + 1) \cdot (\lceil c \log_2(\varepsilon^{-1}) \rceil)$ for node i , and an overall bitstring of
 455 length

$$456 \quad (3.9) \quad \sum_{i=1}^{\tilde{N}} (n(i) + 1) \cdot (\lceil c \log_2(\varepsilon^{-1}) \rceil) \leq 3\tilde{M} \cdot \lceil c \log_2(\varepsilon^{-1}) \rceil$$

457 representing the weights of the graph associated with the network Φ .

458 With (3.8) this shows that the overall number of bits needed to encode the network topol-
 459 ogy and weights is no more than

$$460 \quad (3.10) \quad 5\tilde{M} + 4\tilde{M} \log_2(\tilde{M}) + 2 \log_2(\tilde{M}) + 1 + 3\tilde{M} \cdot \lceil c \log_2(\varepsilon^{-1}) \rceil.$$

462 The network can be recovered by sequentially reading out M, L, d , the \tilde{N}_ℓ , the topology, and
 463 the quantized weights from the overall bitstring. It is not difficult to verify that the individual
 464 steps in the encoding procedure were crafted such that this yields unique recovery. As (3.10)
 465 can be upper-bounded by

$$466 \quad (3.11) \quad C_0 M \log_2(M) \log_2(\varepsilon^{-1})$$

467 for a constant $C_0 > 0$ depending on c and d only, we have constructed an encoder-decoder
 468 pair $(E, D) \in \mathfrak{E}^{\ell(\varepsilon)} \times \mathfrak{D}^{\ell(\varepsilon)}$ with $\ell(\varepsilon)$ satisfying (3.3). This concludes the proof. \square

469 Proposition 3.6 applies to networks that have each weight represented by a finite number of
 470 bits scaling according to $\log_2(\varepsilon^{-1})$ while guaranteeing that the underlying encoder-decoder
 471 pair achieves uniform error ε over \mathcal{C} . We next show that such a compatibility is possible for
 472 networks with activation functions that are either Lipschitz or differentiable such that ρ' is
 473 dominated by a polynomial.

LEMMA 3.7. Let $d, L, k, M \in \mathbb{N}$, $\eta \in (0, 1/2)$, $\Omega \subset \mathbb{R}^d$ be bounded, and let $\rho : \mathbb{R} \rightarrow \mathbb{R}$ be either Lipschitz-continuous or differentiable such that ρ' is dominated by a polynomial. Let $\Phi \in \mathcal{NN}_{L,M,d,\rho}$ with $M \leq \eta^{-k}$ and all its weights be bounded (in absolute value) by η^{-k} . Then, there exist $m \in \mathbb{N}$, depending on k, L , and ρ only, and $\tilde{\Phi} \in \mathcal{NN}_{L,M,d,\rho}$ such that

$$\left\| \tilde{\Phi} - \Phi \right\|_{L^\infty(\Omega)} \leq \eta$$

474 and all weights of $\tilde{\Phi}$ are elements of $\eta^m \mathbb{Z} \cap [-\eta^{-k}, \eta^{-k}]$.

475 *Proof.* We prove the statement for Lipschitz-continuous ρ only. The argument for differ-
476 entiable activation functions with first derivative not growing faster than polynomial is along
477 similar lines.

478 Without loss of generality, we can take the number of nonzero node weights of Φ to be
479 upper-bounded by twice the number of nonzero edge weights. This assumption is justified as
480 else there would be nodes that are not connected to the next layer through an edge of nonzero
481 weight and we could replace the corresponding node weights by zero without altering the
482 mapping Φ . This reduction would then lead to the assumption being justified.

Let $m \in \mathbb{N}$, to be specified later, and denote by $\tilde{\Phi}$ the network that results by replacing
all weights of Φ by a closest element in $\eta^m \mathbb{Z} \cap [-\eta^{-k}, \eta^{-k}]$. Set $C_{\max} := \eta^{-k}$ and denote
the maximum of 1 and the total number of nonzero edge weights plus nonzero node weights
of Φ by C_W . Note that $C_W \leq 3M \leq 3\eta^{-k}$, where the latter inequality is by assumption. For
 $\ell = 1, \dots, L-1$, define $\Phi^\ell : \Omega \rightarrow \mathbb{R}^{N_\ell}$ as

$$\Phi^\ell(x) := \rho(W_\ell \rho(\dots \rho(W_1(x)))) \quad \text{for } x \in \Omega,$$

483 and $\tilde{\Phi}^\ell$ accordingly, and let, for $\ell = 1, \dots, L-1$,

$$484 \quad e_\ell := \left\| \Phi^\ell - \tilde{\Phi}^\ell \right\|_{L^\infty(\Omega, \mathbb{R}^{N_\ell})}, \quad e_L := \left\| \Phi - \tilde{\Phi} \right\|_{L^\infty(\Omega)}.$$

Denote the maximum of 1 and the Lipschitz constant of ρ by C_ρ , set $C_0 := \max\{1, \sup\{|x| : x \in \Omega\}\}$, and let

$$C_\ell := \max \left\{ \left\| \Phi^\ell \right\|_{L^\infty(\Omega, \mathbb{R}^{N_\ell})}, \left\| \tilde{\Phi}^\ell \right\|_{L^\infty(\Omega, \mathbb{R}^{N_\ell})} \right\}, \quad \text{for } \ell = 1, \dots, L-1.$$

486 Then, straightforward, albeit somewhat tedious algebraic manipulations, show that for all
487 $\ell = 2, \dots, L-1$,

$$488 \quad (3.12) \quad e_1 \leq C_0 C_\rho C_W \eta^m, \text{ and } e_\ell \leq C_\rho C_W C_{\ell-1} \eta^m + C_\rho C_W C_{\max} e_{\ell-1}.$$

490 Additionally, we observe that

$$491 \quad (3.13) \quad e_L \leq C_W C_{L-1} \eta^m + C_W C_{\max} e_{L-1}.$$

493 We now bound the quantity C_ℓ for $\ell = 1, \dots, L-1$. A simple computation, exploiting the
494 Lipschitz-continuity of ρ , yields

$$495 \quad C_\ell \leq (|\rho(0)| + C_\rho C_W C_{\max} C_{\ell-1}), \quad \text{for all } \ell = 1, \dots, L-1.$$

Since ρ is continuous on \mathbb{R} we have $|\rho(0)| < \infty$ and thus, by $C_\rho, C_W, C_{\max} \geq 1$, there exists
 $C' > 0$ such that

$$C_\ell \leq C' C_0 (C_\rho C_W C_{\max})^\ell, \quad \text{for all } \ell = 1, \dots, L-1.$$

497 As C_W and C_{\max} are both bounded by η^{-k-2} , it follows that C_ℓ is bounded by η^{-p} for a
 498 $p \in \mathbb{N}$. We can therefore find $n \in \mathbb{N}$ such that for all $\ell = 1, \dots, L-1$,

$$499 \quad (3.14) \quad \max\{C_0 C_\rho C_W, C_W C_{\max}, C_W C_{L-1}, C_\rho C_W C_{\ell-1}, C_\rho C_W C_{\max}\} \leq \frac{\eta^{-n}}{2}.$$

500 Invoking (3.12), we conclude that

$$501 \quad (3.15) \quad e_\ell \leq \frac{\eta^{-n}}{2}(\eta^m + e_{\ell-1}), \quad \text{for all } \ell = 1, \dots, L-1,$$

503 where we set $e_0 = 0$. We proceed by induction to prove that there exists $r \in \mathbb{N}$ such that for
 504 all $\ell = 1, \dots, L-1$,

$$505 \quad (3.16) \quad e_\ell \leq \eta^{m-(\ell-1)n-r}.$$

507 Clearly there exists $r \in \mathbb{N}$ such that $e_1 \leq \eta^{m-r}$. Moreover, one easily verifies that the
 508 existence of an $r \in \mathbb{N}$ such that (3.16) is satisfied for an $\ell \in \{1, \dots, L-2\}$, thanks to (3.15),
 509 implies the existence of an $r \in \mathbb{N}$ such that (3.16) is satisfied for ℓ replaced by $\ell+1$. This
 510 concludes the induction argument.

Using (3.14) and (3.16) in (3.13), we finally obtain

$$e_L \leq \frac{\eta^{m-n}}{2} + \frac{\eta^{m-(L-1)n-r}}{2},$$

511 which yields $e_L \leq \eta$ for sufficiently large m . \square

512 *Remark 3.8.* Note that the weights of the network being elements of $\eta^m \mathbb{Z} \cap [-\eta^{-k}, \eta^{-k}]$
 513 implies that each weight can be represented by no more than $\lceil c \log_2(\eta^{-1}) \rceil$ bits, for some
 514 constant $c > 0$.

515 Proposition 3.6 not only says that the connectivity growth rate of networks achieving
 516 uniform approximation error ε over a function class \mathcal{C} must exceed $\mathcal{O}(\varepsilon^{-1/\gamma^*(\mathcal{C})})$, $\varepsilon \rightarrow 0$,
 517 but its proof, by virtue of constructing an encoder-decoder pair that achieves this growth
 518 rate also provides an achievability result. We next establish a matching strong converse in the
 519 sense of showing that for $\gamma > \gamma^*(\mathcal{C})$, the uniform approximation error remains bounded away
 520 from zero for infinitely many $M \in \mathbb{N}$. To simplify terminology in the sequel, we introduce
 521 the notion of a polynomially bounded variable.

522 **DEFINITION 3.9.** A real variable X depending on the variables $z_i \in D_i \subset \mathbb{R}$, $i =$
 523 $1, \dots, N$, is said to be polynomially bounded in z_1, \dots, z_N , if there exists an N -variate
 524 polynomial π such that $|X| \leq |\pi(z_1, \dots, z_N)|$, for all $z_i \in D_i$, $i = 1, \dots, N$. A set of
 525 real variables $(X_j)_{j \in J}$, each depending on $z_i \in D_i \subset \mathbb{R}$, $i = 1, \dots, N$, is uniformly
 526 polynomially bounded in z_1, \dots, z_N , if there exists an N -variate polynomial π such that
 527 $|X_j| \leq |\pi(z_1, \dots, z_N)|$, for all $j \in J$ and all $z_i \in D_i$, $i = 1, \dots, N$.

528 We will refrain from explicitly specifying the D_i in Definition 3.9 whenever they are clear
 529 from the context.

530 *Remark 3.10.* If $D_i = \mathbb{R} \setminus [-B_i, B_i]$ for some $B_i \geq 1$, $i = 1, \dots, N$, then a variable X
 531 depending on $z_i \in D_i$, $i = 1, \dots, N$, is polynomially bounded in z_1, \dots, z_N if and only if
 532 there exists a $k \in \mathbb{N}$ such that $|X| \leq |z_1^k \cdot z_2^k \cdot \dots \cdot z_N^k|$, for all $z_i \in D_i$.

533 **PROPOSITION 3.11.** Let $d, L \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$ be bounded, π a polynomial, $\mathcal{C} \subset L^2(\Omega)$,
 534 $\rho : \mathbb{R} \rightarrow \mathbb{R}$ either Lipschitz-continuous or differentiable such that ρ' is dominated by a
 535 polynomial. Then, for all $C > 0$ and $\gamma > \gamma^*(\mathcal{C})$, we have that

$$536 \quad (3.17) \quad \sup_{f \in \mathcal{C}} \inf_{\Phi \in \mathcal{NN}_{L, M, d, \rho}^T} \|f - \Phi\|_{L^2(\Omega)} \geq CM^{-\gamma}, \quad \text{for infinitely many } M \in \mathbb{N}.$$

537

Proof. Let $\gamma > \gamma^*(\mathcal{C})$. Assume, towards a contradiction, that (3.17) holds only for finitely many $M \in \mathbb{N}$. Then, there exists a constant $C > 0$ such that the inequality in (3.17) holds for no $M \in \mathbb{N}$ and hence there exists $C' > 0$ so that

$$\sup_{f \in \mathcal{C}} \inf_{\Phi \in \mathcal{NN}_{L,M,d,\rho}^\pi} \|f - \Phi\|_{L^2(\Omega)} \leq C' M^{-\gamma}, \quad \text{for all } M \in \mathbb{N}.$$

Setting $M_\varepsilon := \lceil (\varepsilon/(3C'))^{-1/\gamma} \rceil$, it follows that, for every $f \in \mathcal{C}$ and every $\varepsilon \in (0, 1/2)$, there exists a neural network $\Phi_{\varepsilon,f} \in \mathcal{NN}_{L,M_\varepsilon,d,\rho}^\pi$ such that

$$\|f - \Phi_{\varepsilon,f}\|_{L^2(\Omega)} \leq 2 \sup_{f \in \mathcal{C}} \inf_{\Phi \in \mathcal{NN}_{L,M_\varepsilon,d,\rho}^\pi} \|f - \Phi\|_{L^2(\Omega)} \leq 2C' M_\varepsilon^{-\gamma} \leq \frac{2\varepsilon}{3}.$$

As the weights of $\Phi_{\varepsilon,f}$ are polynomially bounded in M_ε , they are polynomially bounded in ε^{-1} . By Lemma 3.7 and Remark 3.10, there hence exists a network $\tilde{\Phi}_{\varepsilon,f}$ whose weights are represented by no more than $\lceil c \log_2(\varepsilon^{-1}) \rceil$ bits, for some constant $c > 0$, satisfying

$$\|\Phi_{\varepsilon,f} - \tilde{\Phi}_{\varepsilon,f}\|_{L^2(\Omega)} \leq \frac{\varepsilon}{3}.$$

538 Defining

$$539 \quad \mathbf{Learn} : \left(0, \frac{1}{2}\right) \times \mathcal{C} \rightarrow \mathcal{NN}_{\infty,\infty,d,\rho}, \quad (\varepsilon, f) \mapsto \tilde{\Phi}_{\varepsilon,f},$$

it follows that

$$\sup_{f \in \mathcal{C}} \|f - \mathbf{Learn}(\varepsilon, f)\|_{L^2(\Omega)} \leq \varepsilon \quad \text{with} \quad \mathcal{M}(\mathbf{Learn}(\varepsilon, f)) \leq M_\varepsilon \in \mathcal{O}(\varepsilon^{-\frac{1}{\gamma}}), \quad \varepsilon \rightarrow 0.$$

540 The proof is concluded by noting that **Learn** violates Proposition 3.6. \square

541 We can now proceed to the proof of Theorem 3.4.

542 *Proof of Theorem 3.4.* Suppose towards a contradiction that $\gamma_{\mathcal{NN}}^{*,\text{eff}}(\mathcal{C}, \rho) > \gamma^*(\mathcal{C})$. Let
543 $\gamma \in (\gamma^*(\mathcal{C}), \gamma_{\mathcal{NN}}^{*,\text{eff}}(\mathcal{C}, \rho))$. Then, Definition 2.3 implies that there exist a polynomial $\pi, L \in$
544 \mathbb{N} , and $C > 0$ such that

$$545 \quad \sup_{f \in \mathcal{C}} \inf_{\Phi_M \in \mathcal{NN}_{L,M,d,\rho}^\pi} \|f - \Phi_M\|_{L^2(\Omega)} \leq CM^{-\gamma}, \quad \text{for all } M \in \mathbb{N}.$$

547 This, however, constitutes a contradiction to Proposition 3.11. \square

548 We conclude this section with a discussion of the conceptual implications of the re-
549 sults established above. Proposition 3.6 combined with Lemma 3.7 establishes that neural
550 networks with weights polynomially bounded in ε^{-1} and achieving uniform approxima-
551 tion error ε over \mathcal{C} cannot exhibit edge growth rate smaller than $\mathcal{O}(\varepsilon^{-1/\gamma^*(\mathcal{C})})$, $\varepsilon \rightarrow 0$; in
552 other words, a decay of the uniform approximation error, as a function of M , faster than
553 $\mathcal{O}(M^{-\gamma^*(\mathcal{C})})$, $M \rightarrow \infty$, is not possible.

554 Note that requiring uniform approximation error ε only (without imposing the constraint
555 of the network's weights being polynomially bounded in ε^{-1}) can lead to arbitrarily large rate
556 γ as exemplified by Theorem 2.2, which proves the existence of networks realizing an arbi-
557 trarily small approximation error over $L^2([0, 1]^d)$ with a finite number of nodes; in particular,
558 the number of nodes remains constant as $\varepsilon \rightarrow 0$. However, as argued right after Theorem 2.2,
559 these networks necessarily lead to weights that are not polynomially bounded in ε^{-1} .

560 Finally, we remark that the proofs of Theorem 3.4 and Proposition 3.6, by virtue of
 561 explicitly constructing encoder-decoder pairs for neural networks, provide a bound on the
 562 minimax code length of these networks. This, in turn, implies a bound on the networks'
 563 covering numbers, see [26], which, based on classical results from statistical learning theory
 564 (see for example [11]), leads to bounds on the generalization error, see e.g. [3].

565 **4. Transitioning from Representation Systems to Neural Networks.** The remainder
 566 of this paper is devoted to identifying function classes that are optimally representable—
 567 according to Definition 3.5—by neural networks. The mathematical technique we develop in
 568 the process is interesting in its own right as it constitutes a general framework for transferring
 569 results on function approximation through representation systems to results on approximation
 570 by neural networks. In particular, we prove that for a given function class \mathcal{C} and an associated
 571 representation system \mathcal{D} which satisfies certain technical conditions, there exists a neural
 572 network with $\mathcal{O}(M)$ nonzero edge weights that achieves (up to a multiplicative constant) the
 573 same uniform error over \mathcal{C} as a best M -term approximation in \mathcal{D} . This will finally lead to a
 574 characterization of function classes \mathcal{C} that are optimally representable by neural networks in
 575 the sense of Definition 3.5.

576 We start by stating technical conditions on representation systems for the transference
 577 principle outlined above to apply.

DEFINITION 4.1. *Let $d \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$, $\rho : \mathbb{R} \rightarrow \mathbb{R}$, and $\mathcal{D} = (\varphi_i)_{i \in I} \subset L^2(\Omega)$ be a
 representation system. Then, \mathcal{D} is said to be representable by neural networks (with activation
 function ρ), if there exist $L, R \in \mathbb{N}$ such that for all $\eta > 0$ and every $i \in I$, there is a neural
 network $\Phi_{i,\eta} \in \mathcal{NN}_{L,R,d,\rho}$ with*

$$\|\varphi_i - \Phi_{i,\eta}\|_{L^2(\Omega)} \leq \eta.$$

578 *If, in addition, the weights of $\Phi_{i,\eta} \in \mathcal{NN}_{L,R,d,\rho}$ are polynomially bounded in i, η^{-1} , and if
 579 ρ is either Lipschitz-continuous or differentiable such that ρ' is dominated by a polynomial,
 580 then we say that \mathcal{D} is effectively representable by neural networks (with activation function
 581 ρ).*

582 The next result formalizes our transference principle for networks with weights in \mathbb{R} .

THEOREM 4.2. *Let $d \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$, and $\rho : \mathbb{R} \rightarrow \mathbb{R}$. Suppose that $\mathcal{D} = (\varphi_i)_{i \in I} \subset L^2(\Omega)$
 is representable by neural networks. Let $f \in L^2(\Omega)$ and, for $M \in \mathbb{N}$, let $f_M = \sum_{i \in I_M} c_i \varphi_i$,
 $I_M \subset I$, $\#I_M = M$, satisfy*

$$\|f - f_M\|_{L^2(\Omega)} \leq \varepsilon,$$

583 *where $\varepsilon \in (0, 1/2)$. Then, there exist $L \in \mathbb{N}$ (depending on \mathcal{D} only) and a neural network
 584 $\Phi(f, M) \in \mathcal{NN}_{L,M',d,\rho}$ with $M' \in \mathcal{O}(M)$, satisfying*

$$585 \quad (4.1) \quad \|f - \Phi(f, M)\|_{L^2(\Omega)} \leq 2\varepsilon.$$

587 *In particular, for all function classes $\mathcal{C} \subset L^2(\Omega)$, it holds that*

$$588 \quad (4.2) \quad \gamma_{\mathcal{NN}}^*(\mathcal{C}, \rho) \geq \gamma^*(\mathcal{C}, \mathcal{D}).$$

590 *Proof.* By representability of \mathcal{D} according to Definition 4.1, it follows that there exist
 591 $L, R \in \mathbb{N}$, such that for each $i \in I_M$ and for $\eta := \varepsilon / \max\{1, \sum_{i \in I_M} |c_i|\}$, there exists a
 592 neural network $\Phi_{i,\eta} \in \mathcal{NN}_{L,R,d,\rho}$ with

$$593 \quad (4.3) \quad \|\varphi_i - \Phi_{i,\eta}\|_{L^2(\Omega)} \leq \eta.$$

595 Let then $\Phi(f, M)$ be the neural network consisting of the networks $(\Phi_{i,\eta})_{i \in I_M}$ operating in
 596 parallel, all with the same input, and summing their one-dimensional outputs (see Figure 3 in
 597 Section 8 for an illustration) with weights $(c_i)_{i \in I_M}$ according to

$$598 \quad (4.4) \quad \Phi(f, M)(x) := \sum_{i \in I_M} c_i \Phi_{i,\eta}(x), \quad \text{for } x \in \Omega.$$

599 This construction is legitimate as all networks $\Phi_{i,\eta}$ have the same number of layers and
 600 the last layer of a neural network according to Definition 1.1 implements an affine func-
 601 tion only (without subsequent application of the activation function ρ). Then, the fact that
 602 $\Phi(f, M) \in \mathcal{NN}_{L,R,M,d,\rho}$ and application of the triangle inequality together with (4.3) yields
 603 $\|f_M - \Phi(f, M)\|_{L^2(\Omega)} \leq \varepsilon$. Another application of the triangle inequality according to

$$604 \quad \|f - \Phi(f, M)\|_{L^2(\Omega)} \leq \|f - f_M\|_{L^2(\Omega)} + \|f_M - \Phi(f, M)\|_{L^2(\Omega)} \leq 2\varepsilon$$

606 finalizes the proof of (4.1) which by Definitions 1.2 and 1.3 implies (4.2). \square

607 Theorem 4.2 shows that we can restrict ourselves to the approximation of the individual
 608 elements of a representation system by neural networks with the only constraint being that
 609 the number of nonzero edge weights in the individual networks must admit a uniform upper
 610 bound. Theorem 4.2 does, however, not guarantee that the weights of the network $\Phi(f, M)$
 611 can be represented with no more than $\lceil c \log_2(\varepsilon^{-1}) \rceil$ bits when the overall approximation error
 612 is ε . This will again be accomplished through a transfer argument, applied to representation
 613 systems \mathcal{D} satisfying slightly more stringent technical conditions.

THEOREM 4.3. *Let $d \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$ be bounded, and $\mathcal{C} \subset L^2(\Omega)$. Suppose that the
 representation system $\mathcal{D} = (\varphi_i)_{i \in \mathbb{N}} \subset L^2(\Omega)$ is effectively representable by neural networks.
 Then, for all $\gamma < \gamma^{*,\text{eff}}(\mathcal{C}, \mathcal{D})$, there exist a polynomial π , constants $c > 0, L \in \mathbb{N}$, and a map*

$$\mathbf{Learn} : \left(0, \frac{1}{2}\right) \times L^2(\Omega) \rightarrow \mathcal{NN}_{L,\infty,d,\rho}^\pi,$$

614 *such that for every $f \in \mathcal{C}$ the weights in $\mathbf{Learn}(\varepsilon, f)$ can be represented by no more than*
 615 *$\lceil c \log_2(\varepsilon^{-1}) \rceil$ bits while $\|f - \mathbf{Learn}(\varepsilon, f)\|_{L^2(\Omega)} \leq \varepsilon$ and $\mathcal{M}(\mathbf{Learn}(\varepsilon, f)) \in \mathcal{O}(\varepsilon^{-1/\gamma})$,*
 616 *for $\varepsilon \rightarrow 0$.*

617 *Remark 4.4.* Theorem 4.3 implies that if \mathcal{D} optimally represents the function class \mathcal{C}
 618 in the sense of Definition 3.3 and at the same time is effectively representable by neural
 619 networks, then \mathcal{C} is optimally representable by neural networks in the sense of Definition 3.5.

620 *Proof of Theorem 4.3.* Let $M \in \mathbb{N}$ and $\gamma < \gamma^{*,\text{eff}}(\mathcal{C}, \mathcal{D})$. According to Definition 2.1,
 621 there exist constants $C, D > 0$ and a polynomial π such that for every $f \in \mathcal{C}$, there is a subset
 622 $I_M \subset \{1, \dots, \pi(M)\}$, and coefficients $(c_i)_{i \in I_M}$ with $\max_{i \in I_M} |c_i| \leq D$ so that

$$623 \quad (4.5) \quad \left\| f - \sum_{i \in I_M} c_i \varphi_i \right\|_{L^2(\Omega)} \leq \frac{CM^{-\gamma}}{2} =: \frac{\delta_M}{2}.$$

624 We only need to consider the case $\delta_M \leq 1/2$ as will become clear below. By effective
 625 representability according to Definition 4.1, there are $L, R \in \mathbb{N}$ such that for each $i \in I_M$ and
 626 with $\eta := \delta_M / \max\{1, 4 \sum_{i \in I_M} |c_i|\}$, there exists a neural network $\Phi_{i,\eta} \in \mathcal{NN}_{L,R,d,\rho}$ (with
 627 ρ either Lipschitz-continuous or differentiable such that ρ' is dominated by a polynomial)
 628 satisfying

$$629 \quad 630 \quad \|\varphi_i - \Phi_{i,\eta}\|_{L^2(\Omega)} \leq \eta.$$

In addition, the weights of $\Phi_{i,\eta}$ are polynomially bounded in i, η^{-1} . Let then $\Phi(f, M) \in \mathcal{NN}_{L, RM, d, \rho}$ be the neural network consisting of the networks $(\Phi_{i,\eta})_{i \in I_M}$ operating in parallel, according to (4.4). We conclude that

$$\left\| \sum_{i \in I_M} c_i \varphi_i - \Phi(f, M) \right\|_{L^2(\Omega)} \leq \frac{\delta_M}{4}.$$

As the weights of the networks $\Phi_{i,\eta}$ are polynomially bounded in i, η^{-1} and $i \leq \pi(M)$, $\delta_M \sim M^{-\gamma}$, it follows that the weights of $\Phi(f, M)$ are polynomially bounded in δ_M^{-1} ,

$$\left\| \Phi(f, M) - \tilde{\Phi}(f, M) \right\|_{L^2(\Omega)} \leq \frac{\delta_M}{4},$$

631 and all weights of $\tilde{\Phi}(f, M)$ can be represented with no more than $\lceil c \log_2(\delta_M^{-1}) \rceil$ bits, for
632 some $c > 0$. Moreover, we have

$$\begin{aligned} 633 \quad \left\| f - \tilde{\Phi}(f, M) \right\|_{L^2(\Omega)} &\leq \left\| f - \sum_{i \in I_M} c_i \varphi_i \right\|_{L^2(\Omega)} + \left\| \sum_{i \in I_M} c_i \varphi_i - \Phi(f, M) \right\|_{L^2(\Omega)} \\ 634 \quad (4.6) \quad &+ \left\| \Phi(f, M) - \tilde{\Phi}(f, M) \right\|_{L^2(\Omega)} \leq \delta_M = CM^{-\gamma}. \\ 635 \end{aligned}$$

For $\varepsilon \in (0, 1/2)$, we now set

$$\mathbf{Learn}(\varepsilon, f) := \tilde{\Phi}(f, M_\varepsilon),$$

636 where

$$637 \quad (4.7) \quad M_\varepsilon := \left\lceil \left(\frac{C}{\varepsilon} \right)^{\frac{1}{\gamma}} \right\rceil.$$

639 With this choice of M_ε , we have $CM_\varepsilon^{-\gamma} \leq \varepsilon$, which, when used in (4.6), yields

$$640 \quad (4.8) \quad \left\| f - \mathbf{Learn}(\varepsilon, f) \right\|_{L^2(\Omega)} \leq \varepsilon.$$

642 Since, by construction, $\mathbf{Learn}(\varepsilon, f)$ has RM_ε edges and, moreover, $M_\varepsilon \leq C^{1/\gamma} \varepsilon^{-1/\gamma} + 1 \leq$
643 $2C^{1/\gamma} \varepsilon^{-1/\gamma}$, it follows that $\mathbf{Learn}(\varepsilon, f)$ has at most $2RC^{1/\gamma} \varepsilon^{-1/\gamma}$ edges. Moreover, as all
644 weights of $\mathbf{Learn}(\varepsilon, f)$ can be represented by no more than $\lceil c \log_2(\delta_{M_\varepsilon}^{-1}) \rceil$ bits, it follows
645 from $\delta_{M_\varepsilon} \sim M_\varepsilon^{-\gamma} \sim \varepsilon$ that they can be represented by no more than $\lceil c' \log_2(\varepsilon^{-1}) \rceil$ bits, for
646 some $c' > 0$. This concludes the proof. \square

647 **5. All Affine Representation Systems are Effectively Representable by Neural Net-**
648 **works.** This section shows that a large class of representation systems, namely *affine sys-*
649 *tems*, defined below, is effectively representable by neural networks. Affine systems include
650 wavelets, ridgelets, curvelets, shearlets, α -shearlets, and more generally α -molecules. Com-
651 bined with Theorem 4.3 the results in this section establish that any function class that is
652 optimally represented by an arbitrary affine system is optimally represented by neural net-
653 works in the sense of Definition 3.5.

654 Clearly, such strong statements are possible only under restrictions on the choice of the
655 activation function for the approximating neural networks.

656 **5.1. Choice of Activation Function.** We consider two classes of activation functions,
 657 namely sigmoidal functions and smooth approximations of rectified linear units. We start
 658 with the formal definition of sigmoidal activation functions as considered in [12, 43, 45, 7].

659 **DEFINITION 5.1.** A continuous function $\rho : \mathbb{R} \rightarrow \mathbb{R}$ is called a sigmoidal function of
 660 order $k \in \mathbb{N}$, $k \geq 2$, if there exists $C > 0$ such that

$$661 \quad \lim_{x \rightarrow -\infty} \frac{1}{x^k} \rho(x) = 0, \quad \lim_{x \rightarrow \infty} \frac{1}{x^k} \rho(x) = 1, \quad \text{and} \quad |\rho(x)| \leq C(1 + |x|)^k, \quad \text{for } x \in \mathbb{R}.$$

662 A differentiable function ρ is called strongly sigmoidal of order k , if there exist constants
 663 $a, b, C > 0$ such that

$$664 \quad \left| \frac{1}{x^k} \rho(x) \right| \leq C|x|^{-a}, \quad \text{for } x < 0, \quad \left| \frac{1}{x^k} \rho(x) - 1 \right| \leq Cx^{-a}, \quad \text{for } x \geq 0, \quad \text{and}$$

$$665 \quad |\rho(x)| \leq C(1 + |x|)^k, \quad \left| \frac{d}{dx} \rho(x) \right| \leq C|x|^b, \quad \text{for } x \in \mathbb{R}.$$

667 One of the most widely used activation functions is the so-called rectified linear unit (ReLU)
 668 given by $x \mapsto \max\{0, x\}$. The second class of activation functions we consider here are
 669 smooth versions of the ReLU.

670 **DEFINITION 5.2.** Let $\rho : \mathbb{R} \rightarrow \mathbb{R}^+$, $\rho \in C^\infty(\mathbb{R})$, satisfy

$$671 \quad \rho(x) = \begin{cases} 0, & \text{for } x \leq 0, \\ x, & \text{for } x \geq K, \end{cases}$$

672 for some constant $K > 0$. Then, we call ρ an admissible smooth activation function.

673 The reason for considering these two specific classes of activation functions resides in the fact
 674 that neural networks based thereon allow economical representations of multivariate bump
 675 functions, which, in turn, leads to effective representation of all affine systems (built from
 676 bump functions) by neural networks. Approximation of multivariate bump functions using
 677 sparsely connected neural networks is a classical topic in neural network theory [38]. What
 678 is new here is the aspect of quantized weights and rate-distortion optimality.

679 Note that the smooth variants of the ReLU we consider here allow us to build on existing
 680 approximation results for bump functions. We emphasize, however, that smooth functions
 681 can be approximated—in a rate-distortion optimal fashion—by networks based on the regular
 682 ReLU activation function provided that their depth is allowed to scale poly-logarithmically
 683 in the inverse of the approximation error ε [30]. However, in Sections 5.1 and 6 we require
 684 smooth generators of affine systems to be approximated by neural networks with the number
 685 of weights remaining constant as $\varepsilon \rightarrow 0$. This is not possible with the regular ReLU activation
 686 function. We emphasize, though, that all results in this paper apart from those in Sections 5.1
 687 and 6 are also valid for the regular ReLU activation function.

688 A class of bump functions of particular importance in wavelet theory are B -splines. In
 689 [7] it was shown that B -splines can be parsimoniously approximated by neural networks with
 690 sigmoidal activation functions. It is instructive to recall this result. To this end, for $m \in \mathbb{N}$,
 691 we denote the univariate cardinal B -spline of order $m \in \mathbb{N}$ by N_m , i.e., $N_1 = \chi_{[0,1]}$, where
 692 $\chi_{[0,1]}$ denotes the characteristic function of the interval $[0, 1]$, and $N_{m+1} = N_m * \chi_{[0,1]}$,
 693 for all $m \geq 1$. Multivariate B -splines are simply tensor products of univariate B -splines.
 694 Specifically, we denote, for $d \in \mathbb{N}$, the d -dimensional cardinal B -spline of order m by N_m^d .

695 **THEOREM 5.3** ([7], Thm. 4.2). Let $d, m, k \in \mathbb{N}$, and take ρ to be a sigmoidal function
 696 of order $k \geq 2$. Further, let $L := \lceil \log_2(md - d) / \log_2(k) \rceil + 1$. Then, there is $M \in \mathbb{N}$,

697 possibly dependent on d, m, k , such that for all $D, \varepsilon > 0$, there exists a neural network
698 $\Phi_{D,\varepsilon} \in \mathcal{NN}_{L,M,d,\rho}$ with

$$699 \quad \|N_m^d - \Phi_{D,\varepsilon}\|_{L^2([-D,D]^d)} \leq \varepsilon.$$

701 Additionally, we will need to control the weights in the approximating networks $\Phi_{D,\varepsilon}$. We
702 next show that this is, indeed, possible for strongly sigmoidal activation functions.

703 **THEOREM 5.4.** *Let $d, m, k \in \mathbb{N}$, and let ρ be strongly sigmoidal of order $k \geq 2$. Further,
704 let $L := \lceil \log_2(md - d) / \log_2(k) \rceil + 1$. Then, there is an $M \in \mathbb{N}$, possibly dependent on
705 d, m, k , such that for all $D, \varepsilon > 0$, there exists a neural network $\Phi_{D,\varepsilon} \in \mathcal{NN}_{L,M,d,\rho}$ with*

$$706 \quad \|N_m^d - \Phi_{D,\varepsilon}\|_{L^2([-D,D]^d)} \leq \varepsilon.$$

708 *Moreover, the weights of $\Phi_{D,\varepsilon}$ are polynomially bounded in D, ε^{-1} .*

709 *Proof.* The neural network $\Phi_{D,\varepsilon}$ in Theorem 5.3 is explicitly constructed in [7]. Care-
710 fully following the steps in that construction and making explicit use of the strong sigmoidal-
711 ity of ρ , as opposed to plain sigmoidality as in [7], yields the desired result. \square

712 **Remark 5.5.** We observe that the number of edges of the approximating network in The-
713 orem 5.4 does not depend on the approximation error ε .

While Theorem 5.3 demonstrates that a B -spline of order m can be approximated to
arbitrary accuracy by a neural network based on a sigmoidal activation function and of depth
depending on m, d , and the order of sigmoidality of the activation function, we next establish
that for admissible smooth activation functions, exact representation of a general class of
bump functions is possible with a network of three layers only. Before proceeding, we define
for $f \in L^1(\mathbb{R}^d)$, $d \in \mathbb{N}$, the *Fourier transform* of f by

$$\hat{f}(\xi) := \int_{\mathbb{R}^d} f(x) e^{-2\pi i \langle x, \xi \rangle} dx, \text{ for } \xi \in \mathbb{R}^d.$$

714

715 **THEOREM 5.6.** *Let ρ be an admissible smooth activation function. Then, for all $d \in \mathbb{N}$,
716 there exist $M \in \mathbb{N}$ and a neural network $\Phi_\rho \in \mathcal{NN}_{3,M,d,\rho}$ such that*

- 717 (i) Φ_ρ is compactly supported,
- 718 (ii) $\Phi_\rho \in C^\infty(\mathbb{R})$, and
- 719 (iii) $\hat{\Phi}_\rho(\xi) \neq 0$, for all $\xi \in [-3, 3]^d$.

720 *Proof.* We start by constructing an auxiliary function as follows. For $0 < p_1 \leq p_2 \leq p_3$
721 such that $p_1 + p_2 = p_3$, define $t : \mathbb{R} \rightarrow \mathbb{R}$ as

$$722 \quad (5.1) \quad t(x) := \rho(x) - \rho(x - p_1) - \rho(x - p_2) + \rho(x - p_3), \quad x \in \mathbb{R}.$$

724 Then, $t \in C^\infty$ is compactly supported. Letting $q = \|t\|_{L^\infty(\mathbb{R})}$, we define $g : \mathbb{R}^d \rightarrow \mathbb{R}$
725 according to

$$726 \quad (5.2) \quad g(x) := \rho \left(\sum_{i=1}^d t(x_i) - (d-1) \cdot q \right), \quad x \in \mathbb{R}^d.$$

728 By construction, $g \in C^\infty$ is compactly supported. Moreover, g can be realized through a
729 three-layer neural network thanks to its two-step design per (5.1) and (5.2). Since $g \geq 0$ and
730 $g \neq 0$, it follows that $|\hat{g}(0)| > 0$. By continuity of \hat{g} there exists a $\delta > 0$ such that $|\hat{g}(\xi)| > 0$
731 for all $\xi \in [-\delta, \delta]^d$. We now set

$$732 \quad \varphi := g \left(3 \left(\frac{\cdot}{\delta} \right) \right),$$

733 and note that φ can be realized through a three-layer neural network $\Phi_\rho \in \mathcal{NN}_{3,M,d,\rho}$, for
734 some $M \in \mathbb{N}$. As $|\hat{\varphi}(\xi)| > 0$, for all $\xi \in [-3, 3]^d$, Φ_ρ satisfies the desired assumptions. \square

735 **5.2. Invariance to Affine Transformations.** We next leverage Theorems 5.4 and 5.6 to
736 demonstrate that a wide class of representation systems built through affine transformations
737 of B -splines and bump functions as constructed in Theorem 5.6 is effectively representable
738 by neural networks. As a first step towards this general result, we show that representability—
739 in the sense of Definition 4.1—of a single function f by neural networks is invariant to the
740 operation of taking finite linear combinations of affine transformations of f .

741 **PROPOSITION 5.7.** *Let $d \in \mathbb{N}$, $\rho : \mathbb{R} \rightarrow \mathbb{R}$, and $f \in L^2(\mathbb{R}^d)$. Assume that there exist*
742 *$M, L \in \mathbb{N}$ such that for all $D, \varepsilon > 0$, there is $\Phi_{D,\varepsilon} \in \mathcal{NN}_{L,M,d,\rho}$ with*

$$743 \quad (5.3) \quad \|f - \Phi_{D,\varepsilon}\|_{L^2([-D,D]^d)} \leq \varepsilon.$$

744 *Let $A \in \mathbb{R}^{d \times d}$ be full-rank and $b \in \mathbb{R}^d$. Then, there exists $M' \in \mathbb{N}$, depending on M*
745 *and d only, such that for all $E, \eta > 0$, there is $\Psi_{E,\eta} \in \mathcal{NN}_{L,M',d,\rho}$ with*

$$746 \quad \left\| |\det(A)|^{\frac{1}{2}} f(A \cdot - b) - \Psi_{E,\eta} \right\|_{L^2([-E,E]^d)} \leq \eta.$$

747 *Moreover, if the weights of $\Phi_{D,\varepsilon}$ are polynomially bounded in D, ε^{-1} , then the weights of*
748 *$\Psi_{E,\eta}$ are polynomially bounded in $\|A\|_\infty, E, \|b\|_\infty, \eta^{-1}$, where $\|A\|_\infty$ and $\|b\|_\infty$ denote the*
749 *max-norm of A and b , respectively.*

750 *Proof.* By a change of variables, we have for every $\Phi \in \mathcal{NN}_{L,M,d,\rho}$ that

$$751 \quad (5.4) \quad \left\| |\det(A)|^{\frac{1}{2}} (f(A \cdot - b) - \Phi(A \cdot - b)) \right\|_{L^2([-E,E]^d)} = \|f - \Phi\|_{L^2(A \cdot [-E,E]^d - b)},$$

752
753 and there exists a constant M' depending on M and d only such that $|\det(A)|^{1/2} \Phi(A \cdot - b) \in$
754 $\mathcal{NN}_{L,M',d,\rho}$. We furthermore have that

$$755 \quad (5.5) \quad A \cdot [-E, E]^d - b \subset [-(dE\|A\|_\infty + \|b\|_\infty), (dE\|A\|_\infty + \|b\|_\infty)]^d.$$

757 We now set $F = dE\|A\|_\infty + \|b\|_\infty$ and $\Psi_{E,\eta} := |\det(A)|^{1/2} \Phi_{F,\eta}(A \cdot - b)$ and observe that

$$758 \quad \left\| |\det(A)|^{\frac{1}{2}} f(A \cdot - b) - \Psi_{E,\eta} \right\|_{L^2([-E,E]^d)} = \|f - \Phi_{F,\eta}\|_{L^2(A \cdot [-E,E]^d - b)}$$

$$759 \quad \leq \|f - \Phi_{F,\eta}\|_{L^2([-F,F]^d)} \leq \eta,$$

760
761 where we applied the same reasoning as in (5.4) in the first equality and used (5.5) in the
762 first inequality and (5.3) in the second. Moreover, we see that if the weights of $\Phi_{D,\varepsilon}$ are
763 polynomially bounded in D, ε^{-1} , then the weights of $\Psi_{E,\eta}$ are polynomially bounded in
764 $\|A\|_\infty, |\det(A)|, E, \|b\|_\infty, \eta^{-1}$. Since $|\det(A)|$ is polynomially bounded in $\|A\|_\infty$, it follows
765 that the weights of $\Psi_{E,\eta}$ are polynomially bounded in $\|A\|_\infty, E, \|b\|_\infty, \eta^{-1}$. This yields the
766 claim. \square

767 Next, we show that representability by neural networks is preserved under finite linear com-
768 binations of translates.

769 **PROPOSITION 5.8.** *Let $d \in \mathbb{N}$, $\rho : \mathbb{R} \rightarrow \mathbb{R}$, and $f \in L^2(\mathbb{R}^d)$. Assume that there exist*
770 *$M, L \in \mathbb{N}$ such that for all $D, \varepsilon > 0$, there is $\Phi_{D,\varepsilon} \in \mathcal{NN}_{L,M,d,\rho}$ with*

$$771 \quad (5.6) \quad \|f - \Phi_{D,\varepsilon}\|_{L^2([-D,D]^d)} \leq \varepsilon.$$

772 Let $r \in \mathbb{N}$, $(c_i)_{i=1}^r \subset \mathbb{R}$, and $(d_i)_{i=1}^r \subset \mathbb{R}^d$. Then, there exists $M' \in \mathbb{N}$, depending on
 773 M, d , and r only, such that for all $E, \eta > 0$, there is $\Psi_{E,\eta} \in \mathcal{NN}_{L,M',d,\rho}$ with

$$774 \quad (5.7) \quad \left\| \sum_{i=1}^r c_i f(\cdot - d_i) - \Psi_{E,\eta} \right\|_{L^2([-E,E]^d)} \leq \eta.$$

Moreover, if the weights of $\Phi_{D,\varepsilon}$ are polynomially bounded in D, ε^{-1} , then the weights of $\Psi_{E,\eta}$ are polynomially bounded in

$$\sum_{i=1}^r |c_i|, E, \max_{i=1,\dots,r} \|d_i\|_\infty, \eta^{-1}.$$

775 *Proof.* Let $E, \eta > 0$. We start by noting that, for all $D, \varepsilon > 0$,

$$776 \quad \left\| \sum_{i=1}^r c_i f(\cdot - d_i) - \sum_{i=1}^r c_i \Phi_{D,\varepsilon}(\cdot - d_i) \right\|_{L^2([-E,E]^d)} \\ 777 \quad \leq \left(\sum_{i=1}^r |c_i| \right) \cdot \|f - \Phi_{D,\varepsilon}\|_{L^2([-E+d^*], [E+d^*]^d)}, \quad \square \\ 778$$

where $d^* = \max_{i=1,\dots,r} \|d_i\|_\infty$. Setting $D = E + d^*$ and $\varepsilon = \eta / \max\{1, \sum_{i=1}^r |c_i|\}$, and noting that for every $\Phi \in \mathcal{NN}_{L,M,d,\rho}$, the function

$$\Psi := \sum_{i=1}^r c_i \Phi(\cdot - d_i)$$

is in $\mathcal{NN}_{L,M',d,\rho}$ with $M' \in \mathbb{N}$ depending on d, r , and M only, it follows that the network

$$\Psi_{E,\eta} := \sum_{i=1}^r c_i \Phi_{D,\varepsilon}(\cdot - d_i)$$

779 satisfies (5.7). Finally, if the weights of $\Phi_{D,\varepsilon}$ are polynomially bounded in D, ε^{-1} , then the
 780 weights of $\Psi_{E,\eta}$ are polynomially bounded in $\sum_{i=1}^r |c_i|, E, d^*, \eta^{-1}$.

781 Based on the invariance results in Propositions 5.7 and 5.8, we now construct neural networks
 782 which approximate functions with a given number of vanishing moments with arbitrary ac-
 783 curacy. The resulting construction will be crucial in establishing representability of affine
 784 representation systems (see Definition 5.11) by neural networks.

DEFINITION 5.9. Let $R, d \in \mathbb{N}$, and $k \in \{1, \dots, d\}$. A function $g \in C(\mathbb{R}^d)$ is said to possess R directional vanishing moments in the x_k -direction if, for all $\ell \in \{0, \dots, R-1\}$ and all $x_1, \dots, x_{k-1}, x_{k+1}, \dots, x_d \in \mathbb{R}$,

$$\int_{\mathbb{R}} x_k^\ell g(x_1, \dots, x_k, \dots, x_d) dx_k = 0.$$

785 The next result establishes that functions with an arbitrary number of vanishing moments in
 786 a given coordinate direction can be built from suitable linear combinations of translates of a
 787 given continuous function with compact support.

788 LEMMA 5.10. Let $R, d \in \mathbb{N}$, $B > 0$, $k \in \{1, \dots, d\}$, and $f \in C(\mathbb{R}^d)$ with compact
789 support. Then, the function

$$790 \quad (5.8) \quad g(x_1, \dots, x_d) := \sum_{\ell=0}^{R-1} \binom{R-1}{\ell} (-1)^\ell f\left(x_1, \dots, x_k - \frac{\ell}{B}, \dots, x_d\right)$$

791 has R directional vanishing moments in the x_k -direction. Moreover, if $\hat{f}(\xi) \neq 0$ for all
792 $\xi \in [-B, B]^d \setminus \{0\}$, then

$$793 \quad (5.9) \quad \hat{g}(\xi) \neq 0, \quad \text{for all } \xi \in [-B, B]^d \text{ with } \xi_k \neq 0.$$

794 *Proof.* For simplicity of exposition, we consider the case $B = 1$ only. Taking the Fourier
795 transform of (5.8) yields

$$796 \quad (5.10) \quad \hat{g}(\xi) = \sum_{\ell=0}^{R-1} \binom{R-1}{\ell} (-1)^\ell e^{-2\pi i \ell \xi_k} \hat{f}(\xi) = (1 - e^{-2\pi i \xi_k})^{R-1} \cdot \hat{f}(\xi)$$

797

which implies

$$\left(\frac{\partial^\ell}{\partial \xi_k^\ell} \hat{g} \right)_{\xi_k=0} = 0, \quad \text{for all } \ell \in \{0, \dots, R-1\}.$$

798 But, by Definition 5.9, this says precisely that g possesses the desired vanishing moments.
799 Statement (5.9) follows by inspection of (5.10). \square

800 **5.3. Affine Representation Systems.** We are now ready to introduce the general family
801 of representation systems announced earlier in the paper as *affine systems*. This class includes
802 all representation systems based on affine transformations of a given ‘‘mother function’’. Spe-
803 cial cases of affine systems are wavelets, ridgelets, curvelets, shearlets, α -shearlets, and more
804 generally α -molecules, as well as tensor products thereof. The formal definition of affine
805 systems is as follows.

DEFINITION 5.11. Let $d, r, S \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$ be bounded, and let $f \in L^2(\mathbb{R}^d)$ be com-
pactly supported. Let $\delta > 0$, $(c_i^s)_{i=1}^r \subset \mathbb{R}$, for $s = 1, \dots, S$, and let $(d_i)_{i=1}^r \subset \mathbb{R}^d$. Further,
let $A_j \in \mathbb{R}^{d \times d}$, $j \in \mathbb{N}$, be full-rank, with the absolute values of the eigenvalues of A_j bounded
below by 1. Consider the compactly supported functions

$$g_s := \sum_{i=1}^r c_i^s f(\cdot - d_i), \quad s = 1, \dots, S.$$

We define the affine system $\mathcal{D} \subset L^2(\Omega)$ corresponding to $(g_s)_{s=1}^S$ according to

$$\mathcal{D} := \left\{ g_s^{j,b} := \left(|\det(A_j)|^{\frac{1}{2}} g_s(A_j \cdot - \delta b) \right)_{|\Omega} : s = 1, \dots, S, b \in \mathbb{Z}^d, j \in \mathbb{N}, \text{ and } g_s^{j,b} \neq 0 \right\},$$

806 and we refer to f as the generator function of \mathcal{D} .

807 We define the sub-systems $\mathcal{D}_{s,j} := \{g_s^{j,b} \in \mathcal{D} : b \in \mathbb{Z}^d\}$. Since every g_s , $s = 1, \dots, S$,
808 has compact support, $|\mathcal{D}_{s,j}|$ is finite for all $s = 1, \dots, S$ and $j \in \mathbb{N}$. Indeed, we observe that
809 there exists $c_b := c_b((g_s)_{s=1}^S, \delta, d) > 0$ such that for all $s \in \{1, \dots, S\}$, $j \in \mathbb{Z}$, and $b \in \mathbb{Z}^d$,

$$810 \quad (5.11) \quad g_s^{j,b} \in \mathcal{D} \implies \|b\|_\infty \leq c_b \|A_j\|_\infty.$$

811

812 As the $\mathcal{D}_{s,j}$ are finite, we can organize the representation system \mathcal{D} according to

$$813 \quad (5.12) \quad \mathcal{D} = (\varphi_i)_{i \in \mathbb{N}} = (\mathcal{D}_{1,1}, \dots, \mathcal{D}_{S,1}, \mathcal{D}_{1,2}, \dots, \mathcal{D}_{S,2}, \dots),$$

814 where the elements within each sub-system $\mathcal{D}_{s,j}$ may be ordered arbitrarily. This ordering of
815 \mathcal{D} is assumed in the remainder of the paper and will be referred to as *canonical ordering*.

816 Moreover, we note that if there exists $s_o \in \{1, \dots, S\}$ such that g_{s_o} is nonzero, then
817 there is a constant $c_o := c_o((g_s)_{s=1}^S, \delta, d) > 0$ such that

$$818 \quad (5.13) \quad \sum_{s=1}^S |\mathcal{D}_{s,j}| \geq c_o |\det(A_j)|, \text{ for all } j \in \mathbb{N}.$$

820 The next result establishes that all affine systems whose generator functions can be approxi-
821 mated to within arbitrary accuracy by neural networks are (effectively) representable by neu-
822 ral networks.

823 **THEOREM 5.12.** *Let $d \in \mathbb{N}$, $\rho : \mathbb{R} \rightarrow \mathbb{R}$, $\Omega \subset \mathbb{R}^d$ be bounded, and $\mathcal{D} = (\varphi_i)_{i \in \mathbb{N}} \subset$
824 $L^2(\Omega)$ an affine system with generator function f . Suppose that there exist constants $L, R \in$
825 \mathbb{N} such that for all $D, \varepsilon > 0$, there is $\Phi_{D,\varepsilon} \in \mathcal{NN}_{L,R,d,\rho}$ with*

$$826 \quad (5.14) \quad \|f - \Phi_{D,\varepsilon}\|_{L^2([-D,D]^d)} \leq \varepsilon.$$

827 *Then, \mathcal{D} is representable by neural networks with activation function ρ . If, in addition, the*
828 *weights of $\Phi_{D,\varepsilon}$ are polynomially bounded in D, ε^{-1} , and if there exist $a > 0$ and $c > 0$ such*
829 *that*

$$830 \quad (5.15) \quad 1 \geq c \|A_1\|_\infty, \quad \sum_{k=1}^{j-1} |\det(A_k)| \geq c \|A_j\|_\infty^a, \text{ for all } j \in \mathbb{N}, j \geq 2,$$

831 *then \mathcal{D} is effectively representable by neural networks with activation function ρ .*

832 *Proof.* Let $(g_s)_{s=1}^S$ be as in Definition 5.11. If $g_s = 0$ for all $s \in \{1, \dots, S\}$, then $\mathcal{D} =$
833 \emptyset and the result is trivial. Hence, we can assume that there exists at least one $s \in \{1, \dots, S\}$
834 such that $g_s \neq 0$, implying that (5.13) holds.

835 Pick D such that $\Omega \subset [-D, D]^d$. We first show that (5.14) implies representability of \mathcal{D}
836 by neural networks with activation function ρ . To this end, we need to establish the existence
837 of constants $L, R \in \mathbb{N}$ such that for all $i \in \mathbb{N}$ and all $\eta > 0$, there exist $\Phi_{i,\eta} \in \mathcal{NN}_{L,R,d,\rho}$
838 with

$$839 \quad (5.16) \quad \|\varphi_i - \Phi_{i,\eta}\|_{L^2(\Omega)} \leq \eta.$$

841 The elements of \mathcal{D} consist of dilations and translations of f according to

$$842 \quad (5.17) \quad \varphi_i = |\det(A_{j_i})|^{\frac{1}{2}} \left(\sum_{k=1}^r c_k^{s_i} f(A_{j_i} \cdot - \delta b_i - d_k) \right)_{|\Omega},$$

843 for some $r \in \mathbb{N}$ independent of i , and $s_i \in \{1, \dots, S\}$, $j_i \in \mathbb{N}$, and $b_i \in \mathbb{Z}^d$. Thus (5.16)
844 follows directly by Propositions 5.7 and 5.8.

It remains to show that the weights of $\Phi_{D,\varepsilon}$ in (5.14) polynomially bounded in D, ε^{-1}
implies that \mathcal{D} is effectively representable by neural networks with activation function ρ ,
which, by Definition 4.1, means that the weights of $\Phi_{i,\eta}$ are polynomially bounded in i, η^{-1} .
Propositions 5.7 and 5.8 state that the weights of $\Phi_{i,\eta}$ are polynomially bounded in

$$\|A_{j_i}\|_\infty, D, \|b_i\|_\infty, \sum_{k=1}^r |c_k|, \max_{k=1, \dots, r} \|d_k\|_\infty, \eta^{-1}.$$

845 Thanks to (5.11) we have $\|b_i\|_\infty \in \mathcal{O}(\|A_{j_i}\|_\infty)$. Moreover, the quantities D , $\sum_{k=1}^r |c_k|$, and
 846 $\max_{k=1,\dots,r} \|d_k\|_\infty$ do not depend on i . We can thus conclude that the weights of $\Phi_{i,\eta}$ are
 847 polynomially bounded in

$$848 \quad (5.18) \quad \|A_{j_i}\|_\infty, \eta^{-1}.$$

To complete the proof, we need to show that the quantities $\|A_{j_i}\|_\infty$ are polynomially bounded in i . To this end, we first observe that φ_i according to (5.17) satisfies $\varphi_i \in \mathcal{D}_{s_i, j_i}$ for some $s_i \in \{1, \dots, S\}$. Thanks to (5.13) and the canonical ordering (5.12), there exists a constant $c > 0$ such that

$$i \geq c \sum_{k=1}^{j_i-1} |\det(A_k)|, \quad j_i \geq 2.$$

849 We finally appeal to (5.15) to conclude that $\|A_{j_i}\|_\infty$ is polynomially bounded in i for all
 850 $j_i \in \mathbb{N}$, which, together with (5.18), establishes the desired result. \square

851 We remark that condition (5.15) is very weak; in fact, we are not aware of an affine system in
 852 the literature that would violate it.

853 We now proceed to what is probably the central result of this paper, namely that neu-
 854 ral networks provide optimal approximations for all function classes that are optimally ap-
 855 proximated by any affine system with generator function that can be approximated to within
 856 arbitrary accuracy by neural networks.

THEOREM 5.13. *Let $d \in \mathbb{N}$, $\Omega \subset \mathbb{R}^d$ be bounded, $\rho : \mathbb{R} \rightarrow \mathbb{R}$, and $\mathcal{D} = (\varphi_i)_{i \in \mathbb{N}} \subset L^2(\Omega)$ be an affine system with generator function f . Assume that there exist $L, R \in \mathbb{N}$ such that for all $D, \varepsilon > 0$, there is $\Phi_{D,\varepsilon} \in \mathcal{NN}_{L,R,d,\rho}$ satisfying $\|f - \Phi_{D,\varepsilon}\|_{L^2([-D,D]^d)} \leq \varepsilon$. Then, for all function classes $\mathcal{C} \subset L^2(\Omega)$, we have*

$$\gamma_{\mathcal{NN}}^*(\mathcal{C}, \rho) \geq \gamma^*(\mathcal{C}, \mathcal{D}).$$

If, in addition, there is a bivariate polynomial $\tilde{\pi}$ such that the weights of $\Phi_{D,\varepsilon}$ are bounded by $|\tilde{\pi}(D, \varepsilon^{-1})|$, there exist $a > 0$ and $c > 0$ such that (5.15) holds, and \mathcal{C} is optimally represented by \mathcal{D} (according to Definition 3.3), then for all $\gamma < \gamma^(\mathcal{C})$, there exist a constant $c > 0$, a polynomial π , and a map*

$$\mathbf{Learn} : \left(0, \frac{1}{2}\right) \times L^2(\Omega) \rightarrow \mathcal{NN}_{L,\infty,d,\rho}^\pi,$$

857 *such that for every $f \in \mathcal{C}$ the weights in $\mathbf{Learn}(\varepsilon, f)$ can be represented by no more than*
 858 *$\lceil c \log_2(\varepsilon^{-1}) \rceil$ bits while $\|f - \mathbf{Learn}(\varepsilon, f)\|_{L^2(\Omega)} \leq \varepsilon$ and $\mathcal{M}(\mathbf{Learn}(\varepsilon, f)) \in \mathcal{O}(\varepsilon^{-1/\gamma})$,*
 859 *$\varepsilon \rightarrow 0$.*

860 *Proof.* The proof follows directly by combining Theorem 5.12 with Theorems 4.2 and
 861 4.3. \square

862 Theorem 5.13 reveals a remarkable universality and optimality property of neural networks:
 863 All function classes that can be optimally represented by an affine system with generator f
 864 satisfying (5.14) are also optimally representable by neural networks.

865 **6. α -Shearlets and Cartoon-Like Functions.** We next present an explicit pair $(\mathcal{C}, \mathcal{D})$
 866 of function class and representation system satisfying $\gamma_{\mathcal{NN}}^*(\mathcal{C}, \rho) = \gamma^*(\mathcal{C}, \mathcal{D})$. Specifically,
 867 we take α -shearlets as representation system $\mathcal{D} \subset L^2(\mathbb{R}^2)$ and α^{-1} -cartoon-like functions
 868 as function class \mathcal{C} . Cartoon-like functions are piecewise smooth functions with only two
 869 pieces. These pieces are separated by a smooth interface. In a sense, they can be understood

870 as a prototype of a two-dimensional classification function with two homogeneous areas cor-
 871 responding to two classes. Understanding neural network approximation of this function class
 872 is hence relevant to classification tasks in machine learning. We point out that the definition
 873 of α -shearlets in this paper differs slightly from that in [27]. Concretely, relative to [27] our
 874 definition replaces α^{-1} by α so that α -shearlets are a special case of α -molecules, whereas
 875 in [27] α -shearlets are a special case of α^{-1} -molecules. We will need dilation and shearing
 876 matrices defined as

$$877 \quad D_{\alpha,a} := \begin{pmatrix} a & 0 \\ 0 & a^\alpha \end{pmatrix}, \quad J := \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad \text{and} \quad S_k := \begin{pmatrix} 1 & k \\ 0 & 1 \end{pmatrix}.$$

879 This leads us to the following definition which is a slightly modified version of the corre-
 880 sponding definition in [56].

881 **DEFINITION 6.1** ([56]). *For $\delta \in \mathbb{R}^+$, $\alpha \in [0, 1]$, and $f, g \in L^2(\mathbb{R}^2)$, the cone-adapted*
 882 *α -shearlet system $\mathcal{SH}_\alpha(f, g, \delta)$ generated by $f, g \in L^2(\mathbb{R}^2)$ is defined as*

$$883 \quad \mathcal{SH}_\alpha(f, g, \delta) := \mathcal{SH}_\alpha^0(f, g, \delta) \cup \mathcal{SH}_\alpha^1(f, g, \delta),$$

885 where

$$886 \quad \mathcal{SH}_\alpha^0(f, g, \delta) := \{f(\cdot - \delta t) : t \in \mathbb{Z}^2\},$$

$$887 \quad \mathcal{SH}_\alpha^1(f, g, \delta) := \left\{ 2^{\ell \frac{1+\alpha}{2}} g(S_k D_{\alpha, 2^\ell} J^\tau \cdot - \delta t) : \right.$$

$$888 \quad \left. \ell \in \mathbb{N}_0, |k| \leq \lceil 2^{\ell(1-\alpha)} \rceil, t \in \mathbb{Z}^2, k \in \mathbb{Z}, \tau \in \{0, 1\} \right\}.$$

890 Our interest in α -shearlets stems from the fact that they optimally represent α^{-1} -cartoon-like
 891 functions in the sense of Definition 3.3.

892 **DEFINITION 6.2.** *Let $\beta \in [1, 2)$ and $\nu > 0$. Define*

$$893 \quad \mathcal{E}^\beta(\mathbb{R}^2; \nu) = \{f \in L^2(\mathbb{R}^2) : f = f_0 + \chi_B f_1\},$$

895 where $f_0, f_1 \in C^\beta(\mathbb{R}^2)$, $\text{supp } f_0, \text{supp } f_1 \subset (0, 1)^2$, $B \subset [0, 1]^2$, χ_B denotes the char-
 896 acteristic function of B , $\partial B \in C^\beta$, and $\|f_1\|_{C^\beta}, \|f_2\|_{C^\beta}, \|\partial B\|_{C^\beta} < \nu$. The elements of
 897 $\mathcal{E}^\beta(\mathbb{R}^2; \nu)$ are called β -cartoon-like functions.

898 This function class was originally introduced in [20] as a model class for functions gov-
 899 erned by curvilinear discontinuities of prescribed regularity. In this sense, β -cartoon-like
 900 functions provide a convenient model for images governed by edges or for the solutions of
 901 transport equations which often exhibit curvilinear singularities.

902 The optimal exponent $\gamma^*(\mathcal{E}^\beta(\mathbb{R}^2; \nu))$ was found in [20, 28]:

903 **THEOREM 6.3.** *For $\beta \in [1, 2]$ and $\nu > 0$, we have*

$$904 \quad \gamma^*(\mathcal{E}^\beta(\mathbb{R}^2; \nu)) = \frac{\beta}{2}.$$

905 *Proof.* The proof of [20, Theorem 2] demonstrates that a general function class \mathcal{C} has op-
 906 timal exponent $\gamma^*(\mathcal{C}) = (2 - p)/2p$ if \mathcal{C} contains a copy of ℓ_0^p . The result now follows, since
 907 by [28], the function class $\mathcal{E}^\beta(\mathbb{R}^2; \nu)$ does, indeed, contain a copy of ℓ_0^p for $p = 2/(\beta + 1)$. \square

908 Using Proposition 3.6, this result allows us to conclude that neural networks achieving
 909 uniform approximation error ε over the class \mathcal{C} of cartoon-like functions, with weights repre-
 910 sented by no more than $\lceil c \log_2(\varepsilon^{-1}) \rceil$ bits, for some constant $c > 0$, yield an effective best

911 M -edge approximation rate of at most $\beta/2$. Theorem 6.8 below demonstrates achievability
 912 for $\beta = 1/\alpha$, with $\alpha \in [1/2, 1]$.

913 The following theorem states that α -shearlets yield optimal best M -term approximation
 914 rates for α^{-1} -cartoon-like functions.

915 **THEOREM 6.4** ([56], Theorem 6.3 and Remark 6.4). *Let $\alpha \in [1/2, 1]$, $\nu > 0$, $f \in$
 916 $C^{12}(\mathbb{R}^2)$, $g \in C^{32}(\mathbb{R}^2)$, both compactly supported and such that*

- 917 (i) $\widehat{f}(\xi) \neq 0$, for all $|\xi| \leq 1$,
 918 (ii) $\widehat{g}(\xi) \neq 0$, for all $\xi = (\xi_1, \xi_2)^T \in \mathbb{R}^2$ such that $1/3 \leq |\xi_1| \leq 3$ and $|\xi_2| \leq |\xi_1|$,
 (iii) g has at least seven vanishing moments in the x_1 -direction, i.e.,

$$\int_{\mathbb{R}} x_1^\ell g(x_1, x_2) dx_1 = 0, \quad \text{for all } x_2 \in \mathbb{R}, \ell \in \{0, \dots, 6\}.$$

919 Then, there exists $\delta^* > 0$ such that for all $\delta < \delta^*$, the function class $\mathcal{E}^{1/\alpha}(\mathbb{R}^2; \nu)$ is optimally
 920 represented by $\mathcal{SH}_\alpha(f, g, \delta)$.

921 *Remark 6.5.* The assumptions on the smoothness and the number of vanishing moments
 922 of f and g in Theorem 6.4 follow from [56, Eq. 4.9] with $s_1 = 3/2, s_0 = 0, p_0 = q_0 =$
 923 $2/3$, and $|\beta| \leq 4$. While these particular choices allow the statement of the theorem to be
 924 independent of α , it is possible to weaken the assumptions, if a fixed α is considered. For
 925 example, for $\alpha = 1/2$ the smoothness assumptions on f and g reduce to $f \in C^{11}, g \in C^{28}$.

926 As our approximation results for neural networks pertain to bounded domains, we require
 927 a definition of cartoon-like functions on bounded domains.

DEFINITION 6.6. *Let $(0, 1)^2 \subset \Omega \subset \mathbb{R}^2$, $\alpha \in [1/2, 1]$, and $\nu > 0$. We define the set of
 α^{-1} -cartoon-like functions on Ω by*

$$\mathcal{E}^{\frac{1}{\alpha}}(\Omega; \nu) := \left\{ f|_{\Omega} : f \in \mathcal{E}^{\frac{1}{\alpha}}(\mathbb{R}^2; \nu) \right\}.$$

Additionally, for $\delta > 0$, $f, g \in L^2(\mathbb{R}^2)$, we define an α -shearlet system on Ω according to

$$\mathcal{SH}_\alpha(f, g, \delta; \Omega) := \left\{ \phi|_{\Omega} : \phi \in \mathcal{SH}_\alpha(f, g, \delta) \right\}.$$

928 *Remark 6.7.* It is straightforward to check that if $\mathcal{E}^{1/\alpha}(\mathbb{R}^2; \nu)$ is optimally represented
 929 by $\mathcal{SH}_\alpha(f, g, \delta)$, then $\mathcal{E}^{1/\alpha}(\Omega; \nu)$ is optimally represented by $\mathcal{SH}_\alpha(f, g, \delta; \Omega)$.

930 We proceed to the main statement of this section.

931 **THEOREM 6.8.** *Suppose that $(0, 1)^2 \subset \Omega \subset \mathbb{R}^2$ is bounded and $\rho : \mathbb{R} \rightarrow \mathbb{R}$ is either
 932 strongly sigmoidal of order $k \geq 2$ (see Definition 5.1) or an admissible smooth activation
 933 function (see Definition 5.2). Then, for every $\alpha \in [1/2, 1]$, the function class $\mathcal{E}^{1/\alpha}(\Omega; \nu)$ is
 934 optimally representable by a neural network with activation function ρ .*

Proof. Let $\alpha \in [1/2, 1]$ and $\nu > 0$. We first consider the case of ρ strongly sigmoidal
 of order $k \geq 2$. Since the two-dimensional cardinal B -spline of order 34, denoted by N_{34}^2 , is
 32 times continuously differentiable and $\widehat{N}_{34}^2(0) \neq 0$ by construction, we conclude that there
 exists $c > 0$ such that $f := N_{34}^2(c \cdot)$ satisfies $f \in C^{32}(\mathbb{R}^2)$ and $\widehat{f} \neq 0$ for all $\xi \in [-3, 3]^2$.
 Application of Lemma 5.10 then yields the existence of $(c_i)_{i=1}^7 \subset \mathbb{R}$, $(d_i)_{i=1}^7 \subset \mathbb{R}^2$ such
 that $g := \sum_{i=1}^7 c_i f(\cdot - d_i)$ is compactly supported, has seven vanishing moments in the
 x_1 -direction, and $\widehat{g}(\xi) \neq 0$ for all $\xi \in [-3, 3]^2$ such that $\xi_1 \neq 0$. Then, by Theorem 6.4
 and Remark 6.7, there exists $\delta > 0$ such that $\mathcal{SH}_\alpha(f, g, \delta; \Omega)$ is optimal for $\mathcal{E}^{1/\alpha}(\Omega; \nu)$. We
 define

$$\{A_j : j \in \mathbb{N}\} := \left\{ S_k D_{\alpha, 2^\ell} J^\tau : \ell \in \mathbb{N}_0, |k| \leq \lceil 2^{\ell(1-\alpha)} \rceil, \tau \in \{0, 1\} \right\},$$

935 where we order $(A_j)_{j \in \mathbb{N}}$ such that $|\det(A_j)| \leq |\det(A_{j+1})|$, for all $j \in \mathbb{N}$. This construction
 936 implies that the α -shearlet system $\mathcal{SH}_\alpha(f, g, \delta; \Omega)$ is an affine system with generator function
 937 f . Thanks to Theorem 5.4, there exist $L, R \in \mathbb{N}$ such that for all $D, \varepsilon > 0$, there is a network
 938 $\Phi_{D, \varepsilon} \in \mathcal{NN}_{L, R, d, \rho}$ with

$$939 \quad \|f - \Phi_{D, \varepsilon}\|_{L^2([-D, D]^d)} \leq \varepsilon.$$

941 Moreover, the weights of $\Phi_{D, \varepsilon}$ are polynomially bounded in D, ε^{-1} . It is not difficult to
 942 verify that (5.15) holds, and hence Theorem 5.12 yields that $\mathcal{SH}_\alpha(f, g, \delta; \Omega)$ is effectively
 943 representable by neural networks with activation function ρ . Finally, since $\mathcal{E}^{1/\alpha}(\Omega; \nu)$ is
 944 optimally representable by $\mathcal{SH}_\alpha(f, g, \delta; \Omega)$, we conclude with Theorem 4.3 that $\mathcal{E}^{1/\alpha}(\Omega; \nu)$
 945 is optimally representable by neural networks with activation function ρ .

946 It remains to establish the statement for admissible smooth ρ . In this case, by Theo-
 947 rem 5.6 there exist $M \in \mathbb{N}$ and a neural network in $\mathcal{NN}_{3, M, d, \rho}$ which realizes a compactly
 948 supported $f \in C^\infty(\mathbb{R})$ satisfying $\hat{f}(\xi) \neq 0$, for all $\xi \in [-3, 3]^2$. Lemma 5.10 applied to
 949 this f then yields a function g that can be realized by a neural network in $\mathcal{NN}_{3, M', d, \rho}$, for
 950 some $M' \in \mathbb{N}$, has seven vanishing moments in the x_1 -direction, is compactly supported,
 951 and satisfies $g \in C^\infty(\mathbb{R})$, and $\hat{g}(\xi) \neq 0$, for all $\xi \in [-3, 3]^2$ such that $\xi_1 \neq 0$. By Theo-
 952 rem 6.4 and Remark 6.7, there exists $\delta > 0$ such that $\mathcal{E}^{1/\alpha}(\Omega; \nu)$ is optimally representable
 953 by $\mathcal{SH}_\alpha(f, g, \delta; \Omega)$. Note that $\mathcal{SH}_\alpha(f, g, \delta; \Omega)$ is an affine system with generator function
 954 f . Since f can be realized with zero error by a neural network, Theorem 5.12 yields that
 955 $\mathcal{SH}_\alpha(f, g, \delta; \Omega)$ is effectively representable by neural networks with admissible smooth acti-
 956 vation function ρ . Optimality of $\mathcal{SH}_\alpha(f, g, \delta; \Omega)$ for $\mathcal{E}^{1/\alpha}(\Omega; \nu)$ implies, with Theorem 4.3,
 957 that $\mathcal{E}^{1/\alpha}(\Omega; \nu)$ is optimally representable by neural networks with admissible smooth acti-
 958 vation function ρ . \square

959 *Remark 6.9.* Theorem 6.4 requires the generators of the shearlet system guaranteeing
 960 optimal representability of $\mathcal{E}^{1/\alpha}(\Omega; \nu)$, for $1/2 \leq \alpha \leq 1, \nu > 0, \Omega \subset \mathbb{R}^2$ to be very smooth.
 961 On the other hand, Theorem 6.8 demonstrates that optimally-approximating neural networks
 962 are not required to be particularly smooth. Indeed, Theorem 6.8 holds for networks with
 963 differentiable but not necessarily twice differentiable activation function. As the proof of
 964 Theorem 6.8 reveals, such weak assumptions suffice thanks to Theorem 5.4, which demon-
 965 strates that it is possible to approximate arbitrarily smooth B -splines (in L^2 -norm) to within
 966 error ε by neural networks with a number of weights that does not depend on ε as long as the
 967 activation function is strongly sigmoidal.

968 *Remark 6.10.* We observe from the proof of Theorem 6.8 that the depth of the networks
 969 required to achieve optimal approximation depends on the activation function only. Indeed,
 970 for an admissible smooth activation function, inspection of Theorem 5.6 reveals that networks
 971 with three layers can produce optimal approximations in Theorem 6.8. On the other hand,
 972 if a sigmoidal activation function is employed, Theorem 5.4 shows that the construction in
 973 Theorem 6.8 requires a certain minimum depth depending on the order of sigmoidality.

974 **7. Generalization to Manifolds.** Frequently, a function f to be approximated by a
 975 neural network models phenomena on (possibly low-dimensional) immersed submanifolds
 976 $\Gamma \subset \mathbb{R}^d$ of dimension $m < d$. We next briefly outline how our main results can be extended
 977 to this situation. Since analogous results, for the case of wavelets as representation systems,
 978 appear already in [53], we will allow ourselves to be somewhat informal.

Suppose that $f : \Gamma \rightarrow \mathbb{R}$ is compactly supported. Let $(U_i)_{i \in \mathbb{N}} \subset \Gamma$ be an open cover of Γ
 such that for each $i \in \mathbb{N}$ the manifold patch U_i can be parametrized as the graph of a function
 over a subset of the Euclidean coordinates, i.e., there exist coordinates x_{d_1}, \dots, x_{d_m} , open

sets $V_i \subset \mathbb{R}^m$, and smooth mappings

$$\gamma_\ell : \mathbb{R}^m \rightarrow \mathbb{R}, \quad \ell \in \{1, \dots, d\} \setminus \{d_1, \dots, d_m\}$$

979 such that

$$980 \quad U_i = \{\Xi_i(x_{d_1}, \dots, x_{d_m}) := (\gamma_1(x_{d_1}, \dots, x_{d_m}), \dots, x_{d_1}, \dots, \gamma_d(x_{d_1}, \dots, x_{d_m})) : \\ 981 \quad (x_{d_1}, \dots, x_{d_m}) \in V_i\}.$$

983 Take a smooth partition of unity $(h_i)_{i \in \mathbb{N}}$, where $h_i : \Gamma \rightarrow \mathbb{R}$ is smooth with $\text{supp}(h_i) \subset \overline{U_i}$
984 and $\sum_{i \in \mathbb{N}} h_i = 1$. Define the localization of f to U_i by $f_i := fh_i$ such that

$$985 \quad (7.1) \quad f = \sum_{i \in \mathbb{N}} f_i.$$

Every $f_i : U_i \rightarrow \mathbb{R}$ can be reparametrized to

$$\tilde{f}_i : \begin{cases} \mathbb{R}^m & \rightarrow \mathbb{R} \\ (x_{d_1}, \dots, x_{d_m}) & \mapsto f_i \circ \Xi_i(x_{d_1}, \dots, x_{d_m}). \end{cases}$$

986 Suppose that there exist $L, M \in \mathbb{N}$ and neural networks $\tilde{\Phi}_i \in \mathcal{NN}_{L, M, m, \rho}$ such that

$$987 \quad (7.2) \quad \|\tilde{f}_i - \tilde{\Phi}_i\|_{L^2(V_i)} \leq \varepsilon.$$

Then, we can construct a neural network $\Phi_i \in \mathcal{NN}_{L, M+md, d, \rho}$ according to

$$\Phi_i(x) := \tilde{\Phi}_i(P_i x),$$

where P_i denotes the orthogonal projection of x onto the coordinates $(x_{d_1}, \dots, x_{d_m})$. Since P_i is linear, Φ_i is a neural network. Moreover, since P_i is the inverse of the diffeomorphism Ξ_i , we get

$$\|\Phi_i - f_i\|_{L^2(U_i)} \leq C\varepsilon,$$

988 with $C > 0$ depending on the curvature of $\Gamma|_{U_i}$ only. Now we may build a neural network Φ
989 by setting $\Phi := \sum_{i \in \mathbb{N}} \Phi_i$. Combining (7.2) with the observation that, owing to the compact
990 support of f , only a finite number of summands appears in the definition of f , we have
991 constructed a neural network Φ which approximates f on Γ . In summary, we observe the
992 following.

993 *Whenever a function class \mathcal{C} is invariant with respect to diffeomorphisms (in our con-*
994 *struction the functions Ξ_i) and multiplication by smooth functions (in our construction the*
995 *functions h_i), then approximation results on \mathbb{R}^m can be lifted to approximation results on*
996 *m -dimensional submanifolds $\Gamma \subset \mathbb{R}^d$.*

997 Such invariances are, in particular, satisfied for all function classes characterized by a
998 certain smoothness behavior, for example, the class of cartoon-functions as studied in Section
999 6.

1000 **8. Numerical Results.** Our theoretical results show that neural networks realizing uni-
1001 form approximation error ε over a function class $\mathcal{C} \subset L^2(\mathbb{R}^d)$, $d \in \mathbb{N}$, must obey a fundamen-
1002 tal lower bound on the growth rate (as $\varepsilon \rightarrow 0$) of the number of edges of nonzero weight. One
1003 of the most widely used learning algorithms is stochastic gradient descent with the gradient
1004 computed via backpropagation [52]. The purpose of this section is to investigate how this
1005 algorithm fares relative to our lower bound.

1006 Interestingly, our numerical experiments below indicate that for a fixed, sparsely con-
 1007 nected, network topology inspired by the construction of bump functions according to (5.1)
 1008 and (5.2), and with the ReLU as activation function, the stochastic gradient descent algorithm
 1009 generates neural networks that achieve M -edge approximation rates quite close to the fun-
 damental limit. The network topology we prescribe is depicted in Figure 3. The rationale

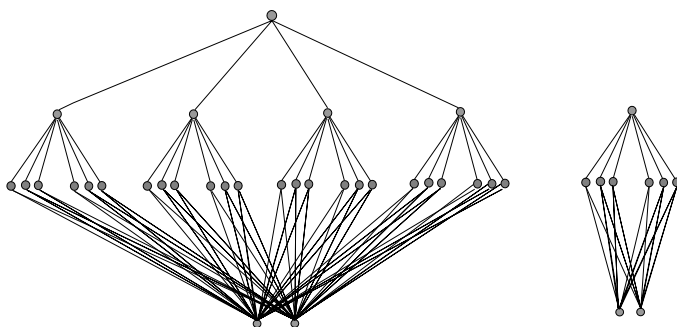


FIG. 3. Left: Topology of the neural network trained using stochastic gradient descent. The network consists of a weighted sum of four subnetworks. Right: A single subnetwork.

1010 for choosing this topology is as follows. As mentioned before, admissible smooth activation
 1011 functions consist of smooth functions which equal a ReLU outside a compact interval. For
 1012 this class of activation functions, the associated α -shearlet generators were constructed from
 1013 a function g as specified in (5.2). Choosing $p_1 = p_2 = 1$ and $p_3 = 2$ in (5.1) yields hat
 1014 functions t . This construction implies that six nodes are required in the first layer in each
 1015 subnetwork. In Figure 3, we see four network realizations of g in parallel. The output layer
 1016 realizes a linear combination of the subnetworks. We now train the network using the stochas-
 1017 tic gradient descent algorithm. Following (5.2) the weights of the second layer remain fixed,
 1018 and the weights in the first and third layers only are trained. Training is performed for two
 1019 different functions, where one is a function with a line singularity (Figure 4(a)), and the other
 1020 one is a cartoon-like function (Figure 5(a)). Specifically, we train the network by drawing
 1021 samples (x_1, x_2) from an equispaced grid in $[-1, 1]^2$. The resulting error is then backprop-
 1022 agated through the network. We repeat this procedure for different network sizes, i.e., for
 1023 different numbers of subnetworks. We start by discussing the results for the function with a
 1024 line singularity depicted in Figure 4(a). The approximation error corresponding to the trained
 1025 neural network is shown in Figure 4(b). The faster than linear decay of the approximation
 1026 error in the semi-logarithmic scale indicates faster than exponential decay with respect to the
 1027 number of edges. This is consistent with the best M -term approximation rate that ridgelets
 1028 yield for piecewise constant functions with line singularities, see [5].

1030 It is interesting to observe that the trained subnetworks yield α -molecules for $\alpha = 0$ (see
 1031 Figures 4(c)-(e)). These functions are constant along one direction and vary along another,
 1032 hence can be considered part of a ridgelet system, which is, in fact, an optimally sparsifying
 1033 representation system for line singularities. Moreover, the orientation of the three learned
 1034 ridge functions matches that of the original function.

1035 In the second experiment, we draw samples from the function depicted in Figure 5(a)
 1036 below, which exhibits a curvilinear singularity. Figures 5(c)-(e) show that the correspond-
 1037 ing trained subnetworks resemble anisotropic molecules with different scales and of different
 1038 orientations. We report, without showing the results, that the decay rate of the corresponding
 1039 approximation error obtained when simply training with different network sizes did not come
 1040 close to the rate of M^{-1} predicted by our theory. However, with a slight adaptation one ob-

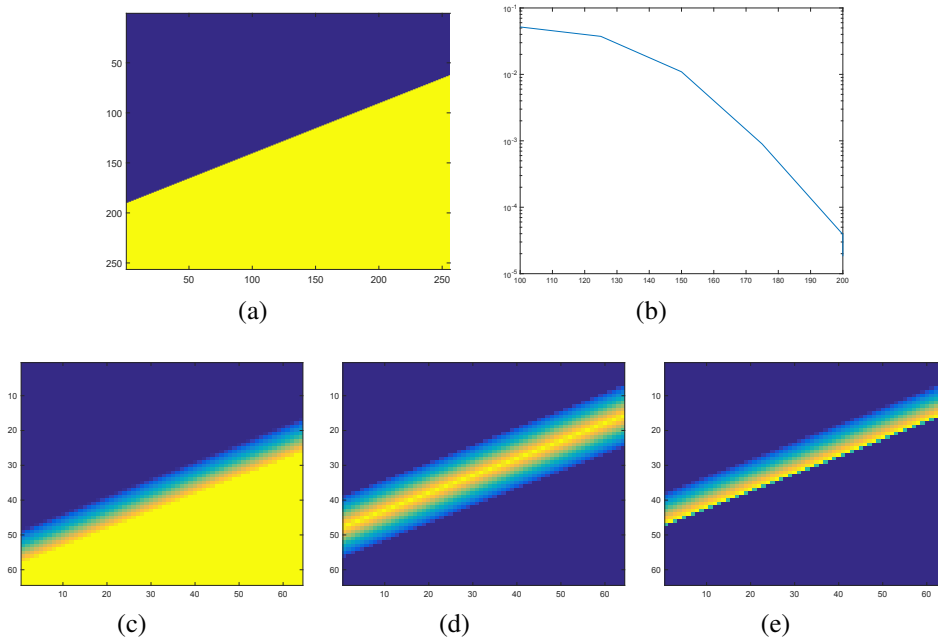


FIG. 4. (a): Function with a line singularity. (b): Approximation error (vertical axis) as a function of the number of edges (horizontal axis). (c)-(e): The functions obtained by restricting to the subnetworks with the largest weights in modulus in the final layer.

1041 tains the result of Figure 5(b), which demonstrates a decay of roughly M^{-1} . The specifics
 1042 of this adaptation are as follows: We first train a large network with ~ 10000 edges, again
 1043 by stochastic gradient descent. Then, the weights in the last layer are optimized using the
 1044 Lasso [55] to obtain a sparse weight vector c^* . We then pick the M largest coefficients of c^*
 1045 and compute the corresponding weighted sum of the associated subnetworks. The resulting
 1046 approximation error is shown in Figure 5(b). Finally, we investigate whether the approxima-
 1047 tion characteristics delivered by this procedure are similar to what would be obtained by best
 1048 M -term approximation with standard shearlet systems. Recall that shearlet elements at high
 1049 scales tend to cluster around singularities [31, 37]. Figures 5(g)-(i) depict the corresponding
 1050 results. Specifically, Figure 5(g) shows the weighted sum of those subnetworks that have the
 1051 largest support. In Figure 5(h), we show weighted sums of subnetworks with medium-sized
 1052 support, and in Figure 5(i) we sum up only the subnetworks with the smallest supports. We
 1053 observe that, indeed, subnetworks of large support approximate the smooth part of the under-
 1054 lying function, whereas the subnetworks associated with small supports resolve the jump
 1055 singularity.

1056 **Acknowledgments.** The authors would like to thank J. Bruna, E. Candès, M. Genzel,
 1057 S. Güntürk, Y. LeCun, K.-R. Müller, H. Rauhut, and F. Voigtländer for interesting discus-
 1058 sions, and D. Perekrestenko and R. Gül for very detailed and insightful comments on the
 1059 manuscript. G. K. and P. P. are grateful to the Faculty of Mathematics at the University of
 1060 Vienna for the hospitality and support during their visits. Moreover, G. K. thanks the De-
 1061 partment of Mathematics at Stanford University whose support allowed for completion of a
 1062 portion of this work. G. K. acknowledges partial support by the Einstein Foundation Berlin,
 1063 the Einstein Center for Mathematics Berlin (ECMath), the European Commission-Project
 1064 DEDALE (contract no. 665044) within the H2020 Framework Program, DFG Grant KU

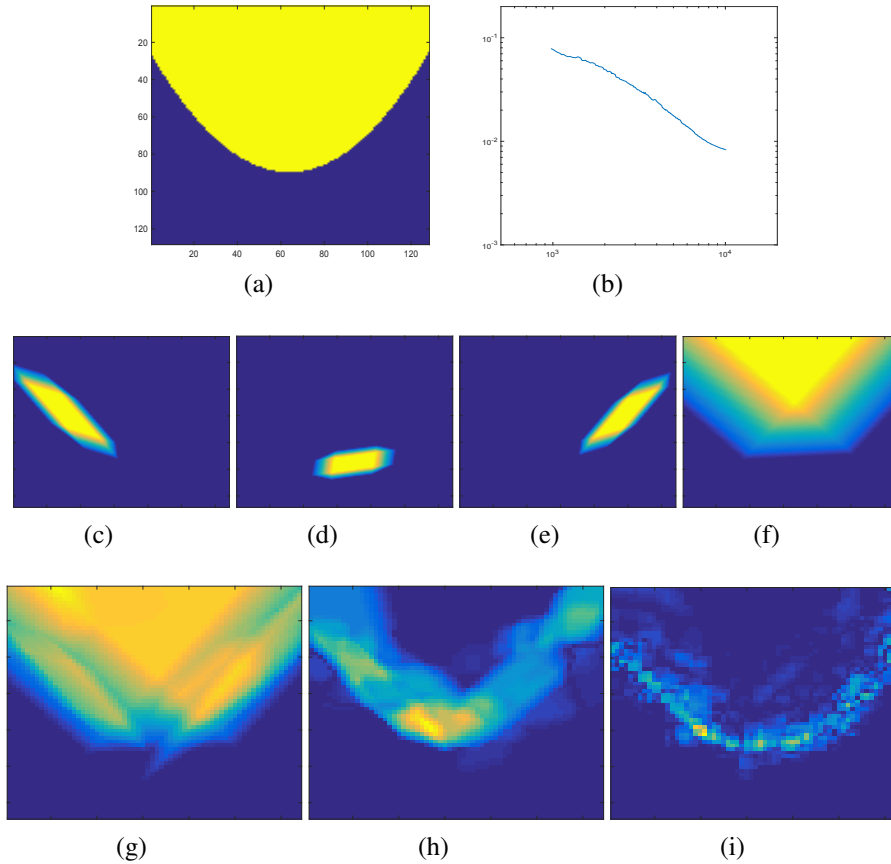


FIG. 5. (a): Function with curvilinear singularity to be approximated by the neural network. (b): Approximation error (vertical axis) as a function of the number of edges (horizontal axis). (c)-(f): Shearlet-like subnetworks. (g): Reconstruction using only the 10 subnetworks whose associated functions have the largest supports. (h): Reconstruction using only subnetworks whose associated functions have medium-sized support. (i): Reconstruction using only subnetworks with associated functions of very small support.

1065 1446/18, DFG-SPP 1798 Grants KU 1446/21 and KU 1446/23, and by the DFG Research
 1066 Center MATHEON “Mathematics for Key Technologies”. G. K. and P. P acknowledge sup-
 1067 port by the DFG Collaborative Research Center TRR 109 “Discretization in Geometry and
 1068 Dynamics”.

1069

REFERENCES

- 1070 [1] A. R. Barron. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Trans. Inf.*
 1071 *Theory*, 39(3):930–945, 1993.
 1072 [2] A. R. Barron. Approximation and estimation bounds for artificial neural networks. *Mach. Learn.*, 14(1):115–
 1073 133, 1994.
 1074 [3] J. Berner, P. Grohs, and A. Jentzen. Analysis of the generalization error: Empirical risk minimization over
 1075 deep artificial neural networks overcomes the curse of dimensionality in the numerical approximation of
 1076 Black-Scholes partial differential equations. *arXiv preprint arXiv:1809.03062*, 2018.
 1077 [4] E. J. Candès. *Ridgelets: Theory and Applications*. Ph.D thesis, Stanford University, 1998.
 1078 [5] E. J. Candès. Ridgelets and the representation of mutilated Sobolev functions. *SIAM J. Math. Anal.*,
 1079 33(2):347–368, 2001.
 1080 [6] E. J. Candès and D. L. Donoho. New tight frames of curvelets and optimal representations of objects with
 1081 piecewise C^2 singularities. *Comm. Pure Appl. Math.*, 57:219–266, 2002.

- 1082 [7] C. K. Chui, X. Li, and H. N. Mhaskar. Neural networks for localized approximation. *Math. Comp.*,
1083 63(208):607–623, 1994.
- 1084 [8] A. Cohen, W. Dahmen, I. Daubechies, and R. A. DeVore. Tree approximation and optimal encoding. *Appl.*
1085 *Comput. Harmon. Anal.*, 11(2):192–226, 2001.
- 1086 [9] N. Cohen, O. Sharir, and A. Shashua. On the expressive power of deep learning: A tensor analysis. In
1087 *Conference on Learning Theory*, pages 698–728, 2016.
- 1088 [10] N. Cohen and A. Shashua. Convolutional rectifier networks as generalized tensor decompositions. In *Inter-*
1089 *national Conference on Machine Learning*, pages 955–963, 2016.
- 1090 [11] F. Cucker and S. Smale. On the mathematical foundations of learning. *Bulletin of the American Mathematical*
1091 *Society*, 39(1):1–49, 2002.
- 1092 [12] G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and*
1093 *Systems*, 2(4):303–314, 1989.
- 1094 [13] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou,
1095 V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T.
1096 Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis. Mastering the game of Go with deep
1097 neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- 1098 [14] I. Daubechies. *Ten Lectures on Wavelets*. SIAM, 1992.
- 1099 [15] L. Demanet and L. Ying. Wave atoms and sparsity of oscillatory patterns. *Appl. Comput. Harmon. Anal.*,
1100 23(3):368–387, 2007.
- 1101 [16] R. A. DeVore. Nonlinear approximation. *Acta Numerica*, 7:51–150, 1998.
- 1102 [17] R. A. DeVore and G. G. Lorentz. *Constructive Approximation*. Springer, 1993.
- 1103 [18] R. A. DeVore, K. Oskolkov, and P. Petrushev. Approximation by feed-forward neural networks. *Ann. Numer.*
1104 *Math.*, 4:261–287, 1996.
- 1105 [19] D. L. Donoho. Unconditional bases are optimal bases for data compression and for statistical estimation.
1106 *Appl. Comput. Harmon. Anal.*, 1(1):100–115, 1993.
- 1107 [20] D. L. Donoho. Sparse components of images and optimal atomic decompositions. *Constr. Approx.*, 17(3):353–
1108 382, 2001.
- 1109 [21] R. Eldan and O. Shamir. The power of depth for feedforward neural networks. In *Proceedings of the 29th*
1110 *Conference on Learning Theory, COLT 2016, New York, USA, June 23-26, 2016*, pages 907–940, 2016.
- 1111 [22] S. Ellacott. Aspects of the numerical analysis of neural networks. *Acta Numer.*, 3:145–202, 1994.
- 1112 [23] K.-I. Funahashi. On the approximate realization of continuous mappings by neural networks. *Neural Net-*
1113 *works*, 2(3):183–192, 1989.
- 1114 [24] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016.
1115 <http://www.deeplearningbook.org>.
- 1116 [25] K. Gröchenig. *Foundations of time-frequency analysis*. Springer Science & Business Media, 2013.
- 1117 [26] P. Grohs. Optimally sparse data representations. In *Harmonic and Applied Analysis*, pages 199–248. Springer,
1118 2015.
- 1119 [27] P. Grohs, S. Keiper, G. Kutyniok, and M. Schäfer. α -molecules. *Appl. Comput. Harmon. Anal.*, 41(1):297–
1120 336, 2016.
- 1121 [28] P. Grohs, S. Keiper, G. Kutyniok, and M. Schäfer. Cartoon approximation with α -curvelets. *J. Fourier Anal.*
1122 *Appl.*, 22(6):1235–1293, 2016.
- 1123 [29] P. Grohs and G. Kutyniok. Parabolic molecules. *Found. Comput. Math.*, 14:299–337, 2014.
- 1124 [30] P. Grohs, D. Perekrestenko, D. Elbrächter, and H. Bölcskei. Deep neural network approximation theory. *IEEE*
1125 *Trans. Inf. Theory*, submitted, 2018.
- 1126 [31] K. Guo, G. Kutyniok, and D. Labate. Sparse multidimensional representations using anisotropic dilation and
1127 shear operators. In *Wavelets and Splines (Athens, GA, 2005)*, pages 189–201. Nashboro Press, Nashville,
1128 TN, 2006.
- 1129 [32] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the*
1130 *IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- 1131 [33] G. Hinton, L. Deng, D. Yu, G. E. Dahl, A. R. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N.
1132 Sainath, and B. Kingsbury. Deep neural networks for acoustic modeling in speech recognition: The
1133 shared views of four research groups. *IEEE Signal Process. Mag.*, 29(6):82–97, 2012.
- 1134 [34] K. Hornik. Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4(2):251 –
1135 257, 1991.
- 1136 [35] K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators.
1137 *Neural Networks*, 2(5):359–366, 1989.
- 1138 [36] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural net-
1139 works. In *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates,
1140 Inc., 2012.
- 1141 [37] G. Kutyniok and W.-Q. Lim. Compactly supported shearlets are optimally sparse. *Journal of Approximation*
1142 *Theory*, 163(11):1564 – 1589, 2011.
- 1143 [38] Y. LeCun. *Modèles connexionnistes de l'apprentissage*. PhD thesis, These de Doctorat, Université Paris 6,

- 1144 1987.
- 1145 [39] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
- 1146 [40] V. Maiorov and A. Pinkus. Lower bounds for approximation by MLP neural networks. *Neurocomputing*,
- 1147 25(1):81–91, 1999.
- 1148 [41] W. McCulloch and W. Pitts. A logical calculus of ideas immanent in nervous activity. *Bull. Math. Biophys.*,
- 1149 5:115–133, 1943.
- 1150 [42] H. Mhaskar and C. Micchelli. Degree of approximation by neural and translation networks with a single
- 1151 hidden layer. *Adv. Appl. Math.*, 16(2):151–183, 1995.
- 1152 [43] H. N. Mhaskar. Approximation properties of a multilayered feedforward artificial neural network. *Advances*
- 1153 *in Computational Mathematics*, 1(1):61–80, Feb 1993.
- 1154 [44] H. N. Mhaskar. Neural networks for optimal approximation of smooth and analytic functions. *Neural Com-*
- 1155 *put.*, 8(1):164–177, 1996.
- 1156 [45] H. N. Mhaskar and C. Micchelli. Approximation by superposition of sigmoidal and radial basis functions.
- 1157 *Adv. Appl. Math.*, 13:350–373, 1992.
- 1158 [46] H. N. Mhaskar and T. Poggio. Deep vs. shallow networks: An approximation theory perspective. *Analysis*
- 1159 *and Applications*, 14(06):829–848, 2016.
- 1160 [47] T. Nguyen-Thien and T. Tran-Cong. Approximation of functions and their derivatives: A neural network
- 1161 implementation with applications. *Appl. Math. Model.*, 23(9):687–704, 1999.
- 1162 [48] E. Ott. *Chaos in dynamical systems*. Cambridge University Press, Cambridge, second edition, 2002.
- 1163 [49] P. Petersen and F. Voigtlaender. Equivalence of approximation by convolutional neural networks and fully-
- 1164 connected networks. arXiv:1809.00973.
- 1165 [50] P. Petersen and F. Voigtlaender. Optimal approximation of piecewise smooth functions using deep ReLU
- 1166 neural networks. *Neural Networks*, 108:296 – 330, 2018.
- 1167 [51] A. Pinkus. Approximation theory of the MLP model in neural networks. *Acta Numer.*, 8:143–195, 1999.
- 1168 [52] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning representations by back-propagating errors.
- 1169 *Nature*, 323(6088):533–536, Oct. 1986.
- 1170 [53] U. Shaham, A. Cloninger, and R. R. Coifman. Provable approximation properties for deep neural networks.
- 1171 *Appl. Comput. Harmon. Anal.*, 44(3):537–557, May 2018.
- 1172 [54] I. Steinwart and A. Christmann. *Support Vector Machines*. Information Science and Statistics. Springer New
- 1173 York, 2008.
- 1174 [55] R. Tibshirani. Regression shrinkage and selection via the Lasso. *J. R. Stat. Soc. Series B Stat. Methodol.*,
- 1175 58(1):267–288, 1996.
- 1176 [56] F. Voigtlaender and A. Pein. Analysis sparsity versus synthesis sparsity for α -shearlets. arXiv:1702.03559.
- 1177 [57] D. Yarotsky. Universal approximations of invariant maps by neural networks. arXiv:1804.10306v1.
- 1178 [58] D. Yarotsky. Error bounds for approximations with deep ReLU networks. *Neural Netw.*, 94:103–114, 2017.
- 1179 [59] D.-X. Zhou. Universality of deep convolutional neural networks. arXiv:1805.10769.