
Nonparametric Classification on Low Dimensional Manifolds using Overparameterized Convolutional Residual Networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Convolutional residual neural networks (ConvResNets), though *overparameter-*
2 *sized*, can achieve remarkable prediction performance in practice, which cannot
3 be well explained by conventional wisdom. To bridge this gap, we study the per-
4 formance of ConvResNeXts, which cover ConvResNets as a special case, trained
5 with weight decay from the perspective of nonparametric classification. Our analy-
6 sis allows for infinitely many building blocks in ConvResNeXts, and shows that
7 weight decay implicitly enforces sparsity on these blocks. Specifically, we consider
8 a smooth target function supported on a low-dimensional manifold, then prove
9 that ConvResNeXts can adapt to the function smoothness and low-dimensional
10 structures and efficiently learn the function without suffering from the curse of
11 dimensionality. Our findings partially justify the advantage of *overparameterized*
12 ConvResNeXts over conventional machine learning models.

13 1 Introduction

14 Deep learning has achieved significant success in various real-world applications. One notable
15 example of this is in the field of image classification, where the winner of the 2017 ImageNet
16 challenge achieved a top-5 error rate of just 2.25% [9] using ConvResNets.

17 Among various deep learning models, ConvResNets have gained widespread popularity in practical
18 applications [2, 8, 20, 28]. Compared to vanilla feedforward neural networks (FNNs), ConvResNets
19 possess two distinct features: convolutional layers and skip connections. Specifically, each block
20 of ConvResNets consists of a subnetwork, called bottleneck, and an identity connection between
21 inconsecutive blocks. The identity connection effectively mitigates the vanishing gradient issue. Each
22 layer of the bottleneck contains several filters (channels) that convolve with the input. Moreover,
23 ConvResNets have various extensions, one of which is ConvResNeXts [25]. This structure generalizes
24 ConvResNets and includes them as a special case. Each building block in ConvResNeXts has a
25 parallel architecture that enables multiple “paths” within the block.

26 There are few theoretical works about ConvResNet, despite its remarkable empirical success. Pre-
27 vious research has focused on the representation power of FNNs [1, 3, 11, 18, 26], while limited
28 literature exists on ConvResNets. Oono and Suzuki [16] developed a representation and statisti-
29 cal estimation theory of ConvResNets, and showed that if the network architecture is appropri-
30 ately designed, ConvResNets with $O(n^{D/(2\alpha+D)})$ blocks can achieve a minimax optimal conver-
31 gence rate $\tilde{O}(n^{-2\alpha/(2\alpha+D)})$ when approximating a C^α function with n samples. Additionally, Liu
32 et al. [14] proved that ConvResNets can universally approximate any function in the Besov space
33 $B_{p,q}^\alpha$ on d -dimensional manifolds with arbitrary accuracy. They improved the convergence rate to
34 $\tilde{O}(n^{-2\alpha/(2\alpha+d)})$ for ConvResNets with $O(n^{d/(2\alpha+d)})$ blocks. Their results only depend on the
35 intrinsic dimension d , rather than the data dimension D .

36 These previous works, however, have limitations in explaining the success of ConvResNets achieved
37 by overparameterization, where the number of blocks can be much larger than the sample size. In

38 practice, the performance of ConvResNets becomes better when they go deeper [8, 24], but the
 39 previous results required a finite number of blocks and thus cannot explain this phenomenon. For
 40 instance, Liu et al. [14] shows that the number of blocks for ConvResNets is $O(n^{d/(2\alpha+d)})$, which is
 41 smaller than the order of the sample size n .

42 To bridge this gap, we study ConvResNeXts under an **overparameterization** regime [25]. We
 43 consider a nonparametric classification problem using ConvResNeXts trained with weight decay. We
 44 prove that even if ConvResNeXts are overparameterized, i.e., the number of blocks is larger than
 45 the order of the sample size n , they can still achieve an asymptotic minimax rate for learning Besov
 46 functions. Specifically, assuming the target function is supported on a d -dimensional manifold and
 47 belongs to the Besov space $B_{p,q}^\alpha$, we prove that the estimator given by the ConvResNeXt class can
 48 converge to the target function at the rate $\tilde{O}(n^{-\frac{\alpha/d}{2\alpha/d+1}(1-o(1))})$ with n samples. Here, weight decay is
 49 a common method in deep learning to reduce overfitting [12, 17]. With this approach, ConvResNeXts
 50 can have infinitely many blocks to achieve arbitrary accuracy, which corresponds to the real-world
 51 applications [8, 24].

52 Our work is partially motivated by Zhang and Wang [27]. However, our work distinguishes itself
 53 through two remarkable technical advancements. Firstly, we develop approximation theory for
 54 ConvResNeXts, while Zhang and Wang [27] only focuses on FNNs. Secondly, we take into account
 55 low-dimensional geometric structures of data. Notably, the statistical rate of convergence in our
 56 theory only depends on the intrinsic dimension d , which circumvents the curse of dimensionality in
 57 Zhang and Wang [27]. Another technical highlight of our paper is bounding the covering number
 58 of weight-decayed ConvResNeXts, which is essential for computing the critical radius of the local
 59 Gaussian complexity. This technique provides a tighter bound than choosing a single radius of the
 60 covering number as in Suzuki [18], Zhang and Wang [27]. To the best of our knowledge, our work is
 61 the first to develop approximation theory and statistical estimation results for ConvResNeXts.

62 2 Architecture of ConvResNeXts

63 In this part, we provide the architecture of ConvResNeXts. This structure has three main features:
 64 residual connections, convolution kernel, and parallel architecture.

65 The building blocks of ConvResNeXts are residual blocks. Given an input \mathbf{x} , each residual block
 66 computes $\mathbf{x} + F(\mathbf{x})$, where F is a subnetwork called bottleneck, consisting of one-sided stride-one
 67 convolutional layers. Figure 2(a) provides a brief illustration of convolution operation $\mathcal{W} \star z$ and its
 68 detailed definition is given in Section A.3 .

69 In ConvResNeXts, a parallel architecture is introduced to each building block, which enables
 70 multiple “paths” in each block. In this paper, we study the ConvResNeXts with rectified linear unit
 71 (ReLU) activation function, i.e., $\text{ReLU}(z) = \max\{z, 0\}$. We next provide the detailed definition of
 72 ConvResNeXts as follows:

73 **Definition 1.** *Let the neural network comprise N residual blocks, each building block has a parallel*
 74 *architecture with M building blocks, and each building block contains L layers. The number of*
 75 *channels is w , and the convolution kernel size is K . Given an input $\mathbf{x} \in \mathbb{R}^D$, a ConvResNeXt with*
 76 *ReLU activation function can be represented as*

$$f(\mathbf{x}) = \mathbf{W}_{\text{out}} \cdot \left(\sum_{m=1}^M f_{N,m} + \text{id} \right) \circ \cdots \circ \left(\sum_{m=1}^M f_{1,m} + \text{id} \right) \circ P(\mathbf{x}),$$

$$f_{n,m} = \mathbf{W}_L^{(n,m)} \star \text{ReLU} \left(\mathbf{W}_{L-1}^{(n,m)} \star \cdots \star \text{ReLU} \left(\mathbf{W}_1^{(n,m)} \star \mathbf{x} \right) \right),$$

77 where id is the identity operator, $P : \mathbb{R}^D \rightarrow \mathbb{R}^{D \times w_0}$ is the padding operator satisfying $P(\mathbf{x}) =$
 78 $[\mathbf{x}, \mathbf{0} \ \dots \ \mathbf{0}] \in \mathbb{R}^{D \times w}$, $\{\mathbf{W}_l^{(n,m)}\}_{l=1}^L$ is a collection of convolution kernels for $n = 1, \dots, N, m =$
 79 $1, \dots, M$, $\mathbf{W}_{\text{out}} \in \mathbb{R}^{w \times L}$ denotes the linear operator for the last layer, and \star is the convolution
 80 operation defined in (6).

81 The structure of ConvResNeXts is shown in Figure 2(b). When $M = 1$, the ConvResNeXt defined in
 82 Definition 1 reduces to a ConvResNet. For the simplicity of notation, we exclude biases in the neural
 83 network structure. This can be compensated by extending the input dimension and padding the input
 84 with a scalar 1 (See Proposition 18 for more details). The channel with 0’s is used to accumulate the
 85 output.

86 3 Theory

87 In this section, we study a binary classification problem on a smooth manifold $\mathcal{M} \subseteq [-1, 1]^D$.
 88 Specifically, we are given i.i.d. samples $\{\mathbf{x}_i, y_i\}_{i=1}^n \sim \mathcal{D}$ where $\mathbf{x}_i \in \mathcal{M}$ and $y_i \in \{0, 1\}$ is the label.
 89 The label y follows the Bernoulli-type distribution

$$\mathbb{P}(y = 1|\mathbf{x}) = \frac{\exp(f^*(\mathbf{x}))}{1 + \exp(f^*(\mathbf{x}))} \quad \text{and} \quad \mathbb{P}(y = 0|\mathbf{x}) = \frac{1}{1 + \exp(f^*(\mathbf{x}))}$$

90 for some $f^* : \mathcal{M} \rightarrow \mathbb{R}$ belonging to the Besov space. Detailed definitions and concepts about smooth
 91 manifold and Besov space are deferred to [Appendix A](#). More specifically, we make the following
 92 assumption on f^* .

93 **Assumption 1.** *Let $0 < p, q \leq \infty$, $d/p < \alpha < \infty$. Assume $f^* \in B_{p,q}^\alpha(\mathcal{M})$ and $\|f^*\|_{B_{p,q}^\alpha(\mathcal{M})} \leq C_F$
 94 for some constant $C_F > 0$.*

95 To learn f^* , we minimize the empirical logistic risk over the training data:

$$\hat{f} = \arg \min_{f \in \mathcal{F}^{\text{Conv}}} \frac{1}{n} \sum_{i=1}^n [y_i \log(1 + \exp(-f(\mathbf{x}_i))) + (1 - y_i) \log(1 + \exp(f(\mathbf{x}_i)))], \quad (1)$$

96 where $\mathcal{F}^{\text{Conv}}$ is some neural network class specified later. For notational simplicity, we denote the
 97 empirical logistic risk function in (1) as $\text{Loss}_n(f)$, and denote the population logistic risk as

$$\mathbb{E}_{\mathcal{D}}[\text{Loss}(f)] = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} y \log(1 + \exp(-f(\mathbf{x}))) + (1 - y) \log(1 + \exp(f(\mathbf{x}))).$$

98 We next specify the class of ConvResNeXts for learning f^* :

$$\mathcal{F}^{\text{Conv}}(N, M, L, K, w, B_{\text{res}}, B_{\text{out}}) = \left\{ f \mid f \text{ is in the form in Definition 1 with } N \text{ residual blocks.} \right.$$

Every residual block has M building blocks with each building block containing L layers.

Each layer has kernel size bounded by K , number of channels bounded by w ,

$$\left. \sum_{n=1}^N \sum_{m=1}^M \sum_{\ell=1}^L \|\mathbf{W}_\ell^{(n,m)}\|_{\mathbb{F}}^2 \leq B_{\text{res}}, \|\mathbf{W}_{\text{out}}\|_{\mathbb{F}}^2 \leq B_{\text{out}}, f(\mathbf{x}) \in [0, 1] \text{ for any } \mathbf{x} \in \mathcal{M}. \right\}$$

99 As can be seen, $\mathcal{F}^{\text{Conv}}$ contains the Frobenius norm constraints of the weights. For the sake of com-
 100 putational convenience in practice, such constraints can be replaced with weight decay regularization
 101 the residual blocks and the last fully connected layer separately. More specifically, we can use the
 102 following alternative formulation:

$$\tilde{f} = \arg \min_{f \in \mathcal{F}^{\text{Conv}}(N, M, L, K, w, \infty, \infty)} \text{Loss}_n(f) + \lambda_1 \sum_{n=1}^N \sum_{m=1}^M \sum_{\ell=1}^L \|\mathbf{W}_\ell^{(n,m)}\|_{\mathbb{F}}^2 + \lambda_2 \|\mathbf{W}_{\text{out}}\|_{\mathbb{F}}^2,$$

103 where $\lambda_1, \lambda_2 > 0$ are properly chosen regularization parameters.

104 3.1 Approximation theory

105 In this section, we provide a universal approximation theory of ConvResNeXts for Besov functions
 106 on a smooth manifold:

107 **Theorem 1.** *For any Besov function f_0 on a smooth manifold satisfying $p, q \geq 1$, $\alpha - d/p > 1$,*

$$\|f_0\|_{B_{p,q}^\alpha(\mathcal{M})} \leq C_F,$$

108 *for any $P > 0$ and any ConvResNeXt class $\mathcal{F}^{\text{Conv}}(N, M, L, K, w, B_{\text{res}}, B_{\text{out}})$ satisfying $L = L' +$
 109 $L_0 - 1, L' \geq 3$, where $L_0 = \lceil \frac{D}{K-1} \rceil$, and*

$$MN \geq C_{\mathcal{M}} P, w \geq C_1(dm + D), B_{\text{res}} \leq C_2 L/K, B_{\text{out}} \leq C_3 C_F^2 ((dm + D)LK)^L (C_{\mathcal{M}} P)^{L-2/p},$$

110 *there exists $f \in \mathcal{F}^{\text{Conv}}(N, M, L, K, w, B_{\text{res}}, B_{\text{out}})$ such that*

$$\|f - f_0\|_{\infty} \leq C_F C_{\mathcal{M}} \left(C_4 P^{-\alpha/d} + C_5 \exp(-C_6 L' \log P) \right), \quad (2)$$

111 *where C_1, C_2, C_3 are universal constants and C_4, C_5, C_6 are constants that only depends on d
 112 and m , d is the intrinsic dimension of the manifold and m is an integer satisfying $0 < \alpha <$
 113 $\min(m, m - 1 + 1/p)$.*

114 The approximation error of the network is bounded by the sum of two terms. The first term is a
 115 polynomial decay term that decreases with the size of the neural network and represents the trailing
 116 term of the B-spline approximation. The second term reflects the approximation error of neural
 117 networks to piecewise polynomials, decreasing exponentially with the number of layers. The proof is
 118 deferred to [Section B.1](#) and the appendix.

119 **3.2 Estimation theory**

120 **Theorem 2.** *Suppose Assumption 1 holds. Set $L = O(\log(n))$ and*

$$MN \geq C_{\mathcal{M}}P, \quad P = O(n^{\frac{1-2/L}{2\alpha/d(1-1/L)+1-2/pL}}).$$

121 *Let \hat{f} be the global minimizer given in (1) with the function class $\mathcal{F} =$*
 122 $\mathcal{F}^{\text{Conv}}(N, M, L, K, w, B_{\text{res}}, B_{\text{out}})$. *Then the estimation error of \hat{f} satisfies*

$$\mathbb{E}_{\mathcal{D}}[\text{Loss}(\hat{f}(\mathbf{x}), y)] \leq \mathbb{E}_{\mathcal{D}}[\text{Loss}(f^*)] + \tilde{O}(n^{-\frac{\alpha/d}{2\alpha/d+1}(1-o(1))}),$$

123 *where $\tilde{O}(\cdot)$ omits the logarithmic term.*

124 The proof for the above theorem is provided in Section B.2. It shows that under weight decay,
 125 the building blocks in a ConvResNeXt are sparse, i.e. only a finite number of blocks contribute
 126 non-trivially to the network even though the model can be overparameterized. This explains why
 127 a ConvResNeXt can generalize well despite overparameterization. Furthermore, we would like to
 128 make the following remarks about the results:

129 • **Strong adaptivity:** By setting the width of the neural network to $w = 2C_1D$, the model can adapt
 130 to any Besov functions on any smooth manifold, provided that $dm \leq D$. This remarkable flexibility
 131 can be achieved simply by tuning the regularization parameter. The cost of overestimating the width
 132 is a slight increase in the estimation error. Considering the immense advantages of this more adaptive
 133 approach, this mild price is well worth paying.

134 • **No curse of dimensionality:** The above error rate only depends polynomially on the ambient
 135 dimension D and exponentially on the intrinsic dimension d . Since in real data, d can be much
 136 smaller than D , this result shows that neural networks can explore the low-dimension structure of
 137 data to overcome the curse of dimensionality.

138 • **Overparameterization is fine:** The number of building blocks in a ConvResNeXt does not
 139 influence the estimation error as long as it is large enough. In other words, our This matches the
 140 empirical observations that neural networks generalize well despite overparameterization.

141 • **Close to minimax rate:** The lower bound of the 1-Lipschitz error for any estimator θ is

$$\min_{\theta} \max_{f^* \in B_{p,q}^{\alpha}} L(\theta(\mathcal{D}), f^*) \gtrsim n^{-\frac{\alpha/d}{2\alpha/d+1}}.$$

142 The proof can be found in Appendix E. Comparing to the minimax rate, we can see that the above
 143 error rate converges to the minimax rate as sample size n grows. In other words, overparameterized
 144 ConvResNeXt can achieve close to the minimax rate in estimating Besov functions. In comparison,
 145 all kernel ridge regression including any NTKs will have a suboptimal rate lower bounded by $\frac{2\alpha-d}{2\alpha}$,
 146 which is suboptimal.

147 **4 Discussion and conclusion**

We compare the Besov space with the Hölder and Sobolev spaces, which are also popular in existing
 literature. The Hölder space $\mathcal{H}^{s,\alpha}$ requires the functions to be differentiable everywhere up to the
 s -th order. The Sobolev space slightly generalizes the Hölder space, but still requires high order
 (weak) differentiability. In contrast, the Besov space $B_{p,q}^s$ does not require weak differentiability,
 and therefore is more general and desirable than the Hölder and Sobolev spaces. Existing work has
 shown that the Besov space can capture important features, such as edges in image processing [10].
 In particular, the Hölder and Sobolev spaces are special cases of the Besov space:

$$\mathcal{H}^{s,\alpha} = W^{s+\alpha,\infty} \subseteq B_{\infty,\infty}^{s+\alpha} \subseteq B_{p,q}^{s+\alpha}$$

148 for any $0 < p, q \leq \infty, s \in \mathbb{N}$ and $\alpha \in (0, 1]$. Due to the generality of the Besov space, existing
 149 literature has been shown that that kernel ridge estimators, including neural tangent kernel only attain
 150 a sub-optimal rate for learning Besov functions [19], which is worse than deep neural networks such
 151 as ConvResNeXts.

152 In this paper, we study the approximation and estimation error of ConvResNeXts. We show that
 153 with proper weight decay, the blocks in a ConvResNeXt converge to a sparse representation, so
 154 the covering number of a ConvResNeXt depends only on the total norm of the parameters and
 155 not on the number of residual blocks, which explains why an overparameterized neural network
 156 generalizes. Assume that the target function is supported on a smooth manifold, the estimation error of
 157 ConvResNeXt depends only weakly (polynomially) on the ambient dimension of the target function.
 158 This result shows that these models do not suffer from the curse of dimensionality, and thus can adapt
 159 to functions on a smooth manifold. While our discussion focuses on binary classification, our result
 160 can be generalized to multi-class classification problems by extending the results to vector-valued
 161 functions.

References

- 162
- 163 [1] Andrew R Barron. Universal approximation bounds for superpositions of a sigmoidal function.
164 *IEEE Transactions on Information theory*, 39(3):930–945, 1993.
- 165 [2] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille.
166 Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and
167 fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):
168 834–848, 2017.
- 169 [3] George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of*
170 *control, signals and systems*, 2(4):303–314, 1989.
- 171 [4] Ronald A DeVore and George G Lorentz. *Constructive approximation*, volume 303. Springer
172 Science & Business Media, 1993.
- 173 [5] Dinh Dũng. Optimal adaptive sampling recovery. *Advances in Computational Mathematics*, 34
174 (1):1–41, 2011.
- 175 [6] Herbert Federer. Curvature measures. *Transactions of the American Mathematical Society*, 93
176 (3):418–491, 1959.
- 177 [7] Daryl Geller and Isaac Z Pesenson. Band-limited localized parseval frames and besov spaces
178 on compact homogeneous manifolds. *Journal of Geometric Analysis*, 21(2):334–371, 2011.
- 179 [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
180 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
181 pages 770–778, 2016.
- 182 [9] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE*
183 *conference on computer vision and pattern recognition*, pages 7132–7141, 2018.
- 184 [10] Stéphane Jaffard, Yves Meyer, and Robert D Ryan. *Wavelets: tools for science and technology*.
185 SIAM, 2001.
- 186 [11] Michael Kohler and Adam Krzyżak. Adaptive regression estimation with multilayer feedforward
187 neural networks. *Nonparametric Statistics*, 17(8):891–913, 2005.
- 188 [12] Anders Krogh and John Hertz. A simple weight decay can improve generalization. *Advances in*
189 *neural information processing systems*, 4, 1991.
- 190 [13] John M Lee. *Riemannian manifolds: an introduction to curvature*, volume 176. Springer
191 Science & Business Media, 2006.
- 192 [14] Hao Liu, Minshuo Chen, Tuo Zhao, and Wenjing Liao. Besov function approximation and
193 binary classification on low-dimensional manifolds using convolutional residual networks. In
194 *International Conference on Machine Learning*, pages 6770–6780. PMLR, 2021.
- 195 [15] Partha Niyogi, Stephen Smale, and Shmuel Weinberger. Finding the homology of submanifolds
196 with high confidence from random samples. *Discrete & Computational Geometry*, 39:419–441,
197 2008.
- 198 [16] Kenta Oono and Taiji Suzuki. Approximation and non-parametric estimation of resnet-type
199 convolutional neural networks. In *International conference on machine learning*, pages 4922–
200 4931. PMLR, 2019.
- 201 [17] Leslie N Smith. A disciplined approach to neural network hyper-parameters: Part 1–learning
202 rate, batch size, momentum, and weight decay. *arXiv preprint arXiv:1803.09820*, 2018.
- 203 [18] Taiji Suzuki. Adaptivity of deep relu network for learning in besov and mixed smooth besov
204 spaces: optimal rate and curse of dimensionality. *arXiv preprint arXiv:1810.08033*, 2018.
- 205 [19] Taiji Suzuki and Atsushi Nitanda. Deep learning is adaptive to intrinsic dimensionality of model
206 smoothness in anisotropic besov space. *Advances in Neural Information Processing Systems*,
207 34:3609–3621, 2021.

- 208 [20] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander Alemi. Inception-v4,
209 inception-resnet and the impact of residual connections on learning. In *Proceedings of the AAAI*
210 *conference on artificial intelligence*, volume 31, 2017.
- 211 [21] H Tribel. Theory of function space ii. *Monographs in Mathematics*, 78, 1992.
- 212 [22] Loring W Tu. Manifolds. In *An Introduction to Manifolds*, pages 47–83. Springer, 2011.
- 213 [23] Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge
214 Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2019. doi:
215 10.1017/9781108627771.013.
- 216 [24] Hongyu Wang, Shuming Ma, Li Dong, Shaohan Huang, Dongdong Zhang, and Furu Wei.
217 Deepnet: Scaling transformers to 1,000 layers. *arXiv preprint arXiv:2203.00555*, 2022.
- 218 [25] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual
219 transformations for deep neural networks. In *Proceedings of the IEEE conference on computer*
220 *vision and pattern recognition*, pages 1492–1500, 2017.
- 221 [26] Dmitry Yarotsky. Error bounds for approximations with deep relu networks. *Neural Networks*,
222 94:103–114, 2017.
- 223 [27] Kaiqi Zhang and Yu-Xiang Wang. Deep learning meets nonparametric regression: Are weight-
224 decayed dnns locally adaptive? *arXiv preprint arXiv:2204.09664*, 2022.
- 225 [28] Qiao Zhang, Zhipeng Cui, Xiaoguang Niu, Shijie Geng, and Yu Qiao. Image segmentation
226 with pyramid dilated convolution based on resnet and u-net. In *Neural Information Processing:*
227 *24th International Conference, ICONIP 2017, Guangzhou, China, November 14-18, 2017,*
228 *Proceedings, Part II 24*, pages 364–372. Springer, 2017.

229 A Background

230 In this section, we introduce some concepts on manifolds. Details can be found in [22] and [13].
 231 Then we provide a detailed definition of the Besov space on smooth manifolds and the convolution
 232 operation.

233 A.1 Smooth manifold

234 Firstly, we briefly introduce manifolds, the partition of unity and reach. Let \mathcal{M} be a d -dimensional
 235 Riemannian manifold isometrically embedded in \mathbb{R}^D with d much smaller than D .

236 **Definition 2** (Chart). *A chart on \mathcal{M} is a pair (U, ϕ) such that $U \subset \mathcal{M}$ is open and $\phi : U \mapsto \mathbb{R}^d$,
 237 where ϕ is a homeomorphism (i.e., bijective, ϕ and ϕ^{-1} are both continuous).*

238 In a chart (U, ϕ) , U is called a coordinate neighborhood, and ϕ is a coordinate system on U .
 239 Essentially, a chart is a local coordinate system on \mathcal{M} . A collection of charts that covers \mathcal{M} is called
 240 an atlas of \mathcal{M} .

241 **Definition 3** (C^k Atlas). *A C^k atlas for \mathcal{M} is a collection of charts $\{(U_i, \phi_i)\}_{i \in \mathcal{A}}$ which satisfies
 242 $\bigcup_{i \in \mathcal{A}} U_i = \mathcal{M}$, and are pairwise C^k compatible:*

$$\phi_i \circ \phi_\beta^{-1} : \phi_\beta(U_i \cap U_\beta) \rightarrow \phi_i(U_i \cap U_\beta) \quad \text{and} \quad \phi_\beta \circ \phi_i^{-1} : \phi_i(U_i \cap U_\beta) \rightarrow \phi_\beta(U_i \cap U_\beta)$$

243 *are both C^k for any $i, \beta \in \mathcal{A}$. An atlas is called finite if it contains finitely many charts.*

244 **Definition 4** (Smooth Manifold). *A smooth manifold is a manifold \mathcal{M} together with a C^∞ atlas.*

245 Classical examples of smooth manifolds are the Euclidean space, the torus, and the unit sphere.
 246 Furthermore, we define C^s functions on a smooth manifold \mathcal{M} as follows:

247 **Definition 5** (C^s functions on \mathcal{M}). *Let \mathcal{M} be a smooth manifold and $f : \mathcal{M} \rightarrow \mathbb{R}$ be a function on
 248 \mathcal{M} . A function $f : \mathcal{M} \rightarrow \mathbb{R}$ is C^s if for any chart (U, ϕ) on \mathcal{M} , the composition $f \circ \phi^{-1} : \phi(U) \rightarrow \mathbb{R}$
 249 is a continuously differentiable up to order s .*

250 We next define the C^∞ partition of unity, which is an important tool for studying functions on
 251 manifolds.

252 **Definition 6** (Partition of Unity, Definition 13.4 in [22]). *A C^∞ partition of unity on a manifold \mathcal{M}
 253 is a collection of C^∞ functions $\{\rho_i\}_{i \in \mathcal{A}}$ with $\rho_i : \mathcal{M} \rightarrow [0, 1]$ such that for any $\mathbf{x} \in \mathcal{M}$,*

254 *1. there is a neighbourhood of \mathbf{x} where only a finite number of the functions in $\{\rho_i\}_{i \in \mathcal{A}}$ are
 255 nonzero;*

256 *2. $\sum_{i \in \mathcal{A}} \rho_i(\mathbf{x}) = 1$.*

257 An open cover of a manifold \mathcal{M} is called locally finite if every $\mathbf{x} \in \mathcal{M}$ has a neighborhood that
 258 intersects with a finite number of sets in the cover. The following proposition shows that a C^∞
 259 partition of unity for a smooth manifold always exists.

260 **Proposition 3** (Existence of a C^∞ partition of unity, Theorem 13.7 in [22]). *Let $\{U_i\}_{i \in \mathcal{A}}$ be a locally
 261 finite cover of a smooth manifold \mathcal{M} . Then there is a C^∞ partition of unity $\{\rho_i\}_{i=1}^\infty$ where every ρ_i
 262 has a compact support such that $\text{supp}(\rho_i) \subset U_i$.*

263 Let $\{(U_i, \phi_i)\}_{i \in \mathcal{A}}$ be a C^∞ atlas of \mathcal{M} . Proposition 3 guarantees the existence of a partition of unity
 264 $\{\rho_i\}_{i \in \mathcal{A}}$ such that ρ_i is supported on U_i . To characterize the curvature of a manifold, we adopt the
 265 geometric concept: reach.

266 **Definition 7** (Reach [6, 15]). *Denote*

$$G = \left\{ \mathbf{x} \in \mathbb{R}^D : \exists \mathbf{p} \neq \mathbf{q} \in \mathcal{M} \text{ such that } \|\mathbf{x} - \mathbf{p}\|_2 = \|\mathbf{x} - \mathbf{q}\|_2 = \inf_{\mathbf{y} \in \mathcal{M}} \|\mathbf{x} - \mathbf{y}\|_2 \right\}.$$

*as the set of points with at least two nearest neighbors on \mathcal{M} . The closure of G is called the medial
 axis of \mathcal{M} . Then the reach of \mathcal{M} is defined as*

$$\tau = \inf_{\mathbf{x} \in \mathcal{M}} \inf_{\mathbf{y} \in G} \|\mathbf{x} - \mathbf{y}\|_2.$$

267 Reach has a simple geometrical interpretation: for every point $\mathbf{x} \in \mathcal{M}$, the osculating circle's radius
 268 is at least τ . A large reach for \mathcal{M} indicates that the manifold changes slowly.

269 **A.2 Besov functions on a smooth manifold**

270 We next define the Besov function space on the smooth manifold \mathcal{M} , which generalizes more
 271 elementary function spaces such as the Sobolev and Hölder spaces.

272 Roughly speaking, functions in the Besov space are only required to have weak derivatives with
 273 bounded total variation. For example, consider a wiggly piecewise linear function as shown in Figure
 274 1. Its derivative can go to infinite while the total variation of the function is upper bounded. Therefore,
 275 according to the definition of the Besov space in Definition 9, the function given in Figure 1 is Besov,
 276 but not Hölder.

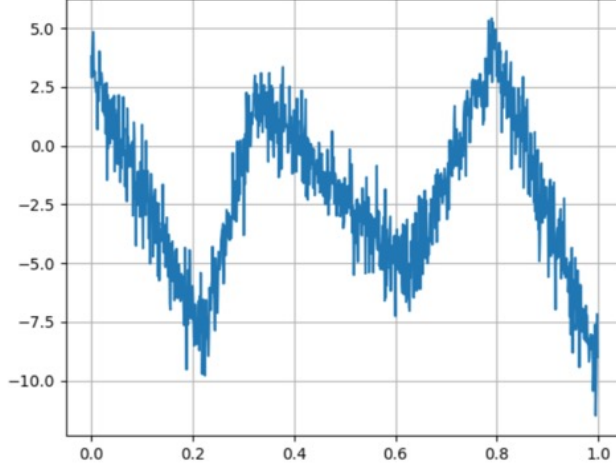


Figure 1: A piecewise linear function which is Besov.

277 To define Besov functions rigorously, we first introduce the modulus of smoothness.

278 **Definition 8** (Modulus of Smoothness [4, 18]). *Let $\Omega \subset \mathbb{R}^D$. For a function $f : \mathbb{R}^D \rightarrow \mathbb{R}$ be in*
 279 *$L^p(\Omega)$ for $p > 0$, the r -th modulus of smoothness of f is defined by*

$$w_{r,p}(f, t) = \sup_{\|\mathbf{h}\|_2 \leq t} \|\Delta_{\mathbf{h}}^r(f)\|_{L^p}, \text{ where}$$

$$\Delta_{\mathbf{h}}^r(f)(\mathbf{x}) = \begin{cases} \sum_{j=0}^r \binom{r}{j} (-1)^{r-j} f(\mathbf{x} + j\mathbf{h}) & \text{if } \mathbf{x} \in \Omega, \mathbf{x} + r\mathbf{h} \in \Omega, \\ 0 & \text{otherwise.} \end{cases}$$

Definition 9 (Besov Space $B_{p,q}^\alpha(\Omega)$). *For $0 < p, q \leq \infty, \alpha > 0, r = \lfloor \alpha \rfloor + 1$, define the seminorm*
 $|\cdot|_{B_{p,q}^\alpha}$ *as*

$$|f|_{B_{p,q}^\alpha(\Omega)} := \begin{cases} \left(\int_0^\infty (t^{-\alpha} w_{r,p}(f, t))^q \frac{dt}{t} \right)^{\frac{1}{q}} & \text{if } q < \infty, \\ \sup_{t>0} t^{-\alpha} w_{r,p}(f, t) & \text{if } q = \infty. \end{cases}$$

280

281 *The norm of the Besov space $B_{p,q}^\alpha(\Omega)$ is defined as $\|f\|_{B_{p,q}^\alpha(\Omega)} := \|f\|_{L^p(\Omega)} + |f|_{B_{p,q}^\alpha(\Omega)}$. Then the*
 282 *Besov space is defined as $B_{p,q}^\alpha(\Omega) = \{f \in L^p(\Omega) \mid \|f\|_{B_{p,q}^\alpha} < \infty\}$.*

283 Moreover, we show that functions in the Besov space can be decomposed using B-spline basis
 284 functions in the following proposition.

285 **Proposition 4** (Decomposition of Besov functions). *Any function f in the Besov space $B_{p,q}^\alpha, \alpha > d/p$*
 286 *can be decomposed using B-spline of order $m, m > \alpha$: for any $\mathbf{x} \in \mathbb{R}^d$, we have*

$$f(\mathbf{x}) = \sum_{k=0}^{\infty} \sum_{\mathbf{s} \in J(k)} c_{k,\mathbf{s}}(f) M_{m,k,\mathbf{s}}(\mathbf{x}), \quad (3)$$

287 where $J(k) := \{2^{-k} \mathbf{s} : \mathbf{s} \in [-m, 2^k + m]^d \subset \mathbb{Z}^d\}$, $M_{m,k,\mathbf{s}}(\mathbf{x}) := M_m(2^k(\mathbf{x} - \mathbf{s}))$, and $M_k(\mathbf{x}) =$
 288 $\prod_{i=1}^d M_k(x_i)$ is the cardinal B-spline basis function which can be expressed as a polynomial:

$$\begin{aligned} M_m(z) &= \frac{1}{m!} \sum_{j=1}^{m+1} (-1)^j \binom{m+1}{j} (z-j)_+^m \\ &= ((m+1)/2)^m \frac{1}{m!} \sum_{j=1}^{m+1} (-1)^j \binom{m+1}{j} \left(\frac{z-j}{(m+1)/2} \right)_+^m. \end{aligned} \quad (4)$$

289 We next define $B_{p,q}^\alpha$ functions on \mathcal{M} .

290 **Definition 10** ($B_{p,q}^\alpha$ Functions on \mathcal{M} [7, 21]). Let \mathcal{M} be a compact smooth manifold of dimension
 291 d . Let $\{(U_i, \phi_i)\}_{i=1}^{C_{\mathcal{M}}}$ be a finite atlas on \mathcal{M} and $\{\rho_i\}_{i=1}^{C_{\mathcal{M}}}$ be a partition of unity on \mathcal{M} such that
 292 $\text{supp}(\rho_i) \subset U_i$. A function $f : \mathcal{M} \rightarrow \mathbb{R}$ is in $B_{p,q}^\alpha(\mathcal{M})$ if

$$\|f\|_{B_{p,q}^\alpha(\mathcal{M})} := \sum_{i=1}^{C_{\mathcal{M}}} \|(f\rho_i) \circ \phi_i^{-1}\|_{B_{p,q}^\alpha(\mathbb{R}^d)} < \infty. \quad (5)$$

293 Since ρ_i is supported on U_i , the function $(f\rho_i) \circ \phi_i^{-1}$ is supported on $\phi(U_i)$. We can extend $(f\rho_i) \circ \phi_i^{-1}$
 294 from $\phi(U_i)$ to \mathbb{R}^d by setting the function to be 0 on $\mathbb{R}^d \setminus \phi(U_i)$. The extended function lies in the
 295 Besov space $B_{p,q}^s(\mathbb{R}^d)$ [21, Chapter 7].

296 A.3 Architecture of ConvResNeXt

297 In this part, we present the formulation of the one-sided stride-one convolution in ConvResNeXts.
 298 Let $\mathcal{W} = \{\mathcal{W}_{j,k,l}\} \in \mathbb{R}^{w' \times K \times w}$ be a convolution kernel with output channel size w' , kernel size K
 299 and input channel size w . For $\mathbf{z} \in \mathbb{R}^{D \times w}$, the convolution of \mathcal{W} with \mathbf{z} gives $\mathbf{y} \in \mathbb{R}^{D \times w'}$ such that

$$\mathbf{y} = \mathcal{W} * \mathbf{z}, \quad y_{i,j} = \sum_{k=1}^K \sum_{l=1}^w \mathcal{W}_{j,k,l} z_{i+k-1,l}, \quad (6)$$

300 where $1 \leq i \leq D, 1 \leq j \leq w'$ and we set $z_{i+k-1,l} = 0$ for $i+k-1 > D$, as demonstrated in
 301 Figure 2(a).

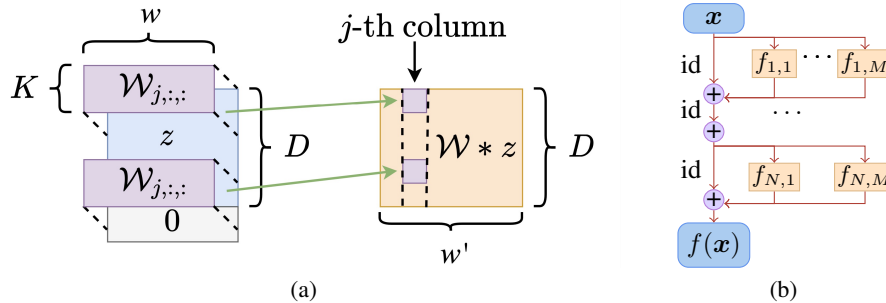


Figure 2: (a) Demonstration of the convolution operation $\mathcal{W} * z$, where the input is $z \in \mathbb{R}^{D \times w}$, and the output is $\mathcal{W} * z \in \mathbb{R}^{D \times w'}$. Here $\mathcal{W}_{j,:}$ is a $D \times w$ matrix for the j -th output channel. (b) Demonstration of the ConvResNeXt. $f_{1,1} \dots f_{N,M}$ are the building blocks, each building block is a convolution neural network.

302 B Proof overview

303 B.1 Approximation error

304 We follow the method in Liu et al. [14] to construct a neural network that achieves the approximation
 305 error we claim. It is divided into the following steps:

- 306 • **Step 1: Decompose the target function into the sum of locally supported functions.**

307 In this work, we adopt a similar approach to [14] and partition \mathcal{M} using a finite number of open
 308 balls on \mathbb{R}^D . Specifically, we define $B(\mathbf{c}_i, r)$ as the set of unit balls with center \mathbf{c}_i and radius r such
 309 that their union covers the manifold of interest, i.e., $\mathcal{M} \subseteq \cup_{i=1}^{C_{\mathcal{M}}} B(\mathbf{c}_i, r)$. This allows us to partition

310 the manifold into subregions $U_i = B(\mathbf{c}_i, r) \cap \mathcal{M}$, and further decompose a smooth function on the
 311 manifold into the sum of locally supported smooth functions with linear projections. The existence of
 312 function decomposition is guaranteed by the existence of partition of unity stated in [Proposition 3](#).
 313 See [Section C.1](#) for the detail.

314 • **Step 2: Locally approximate the decomposed functions using cardinal B-spline basis functions.**
 315 In the second step, we decompose the locally supported Besov functions achieved in the first step
 316 using B-spline basis functions. The existence of the decomposition was proven by Dũng [5], and was
 317 applied in a series of works [27, 18, 14]. The difference between our result and previous work is that
 318 we define a norm on the coefficients and bound this norm, instead of bounding the maximum value.
 319 The detail is deferred to [Section C.2](#).

320 • **Step 3: Approximate the polynomial functions using neural networks.** In this section, we follow
 321 the method in Zhang and Wang [27], Suzuki [18], Liu et al. [14] and show that neural networks can
 322 be used to approximate polynomial functions, including B-spline basis functions and the distance
 323 function. The key technique is to use a neural network to approximate square function and multiply
 324 function [1]. The detail is deferred to the appendix. Specifically, [Lemma 17](#) proves that a neural
 325 network with width $w = O(dm)$ and depth L can approximate B-spline basis functions, and the error
 326 decreases exponentially with L ; Similarly, [Proposition 9](#) shows that a neural network with width
 327 $w = O(D)$ can approximately calculate the distance between two points $d^2(\mathbf{x}; \mathbf{c})$, with precision
 328 decreasing exponentially with the depth.

329 • **Step 4: Use a ConvResNeXt to Approximate the target function.** Using the results above, the
 330 target function can be (approximately) decomposed as

$$\sum_{i=1}^{C_{\mathcal{M}}} \sum_{j=1}^P a_{i,k_j,s_j} M_{m,k_j,s_j} \circ \phi_i \times \mathbf{1}(\mathbf{x} \in B(\mathbf{c}_i, r)). \quad (7)$$

331 We first demonstrate that a ReLU neural network taking two scalars a, b as the input, denoted as
 332 $a \tilde{\times} b$, can approximate

$$y \times \mathbf{1}(\mathbf{x} \in B_{r,i}),$$

333 where $\tilde{\times}$ satisfy that $y \tilde{\times} 1 = y$ for all y , and $y \tilde{\times} \tilde{x} = 0$ if any of x or y is 0, and the soft indicator
 334 function $\tilde{\mathbf{1}}(\mathbf{x} \in B_{r,i})$ satisfy $\tilde{\mathbf{1}}(\mathbf{x} \in B_{r,i}) = 1$ when $x \in B_{r,i}$, and $\tilde{\mathbf{1}}(\mathbf{x} \in B_{r,i}) = 0$ when
 335 $x \notin B_{r+\Delta,i}$. The detail is deferred to [Section C.3](#).

Then, we show that it is possible to construct $MN = C_{\mathcal{M}}P$ number of building blocks, such that
 each building block is a feedforward neural network with width $C_1(md + D)$ and depth L , where m
 is an interger satisfying $0 < \alpha < \min(m, m - 1 + 1/p)$. The k -th building block (the position of
 the block does not matter) approximates

$$a_{i,k_j,s_j} M_{m,k_j,s_j} \circ \phi_i \times \mathbf{1}(\mathbf{x} \in B(\mathbf{c}_i, r)),$$

336 where $i = \text{ceiling}(k/N)$, $j = \text{rem}(k, N)$. Each building block has where a sub-block with width
 337 D and depth $L - 1$ approximates the chart selection, a sub-block with width md and depth $L - 1$
 338 approximates the B-spline function, and the last layer approximates the multiply function. The norm
 339 of this block is bounded by

$$\sum_{\ell=1}^L \|\mathbf{W}_{\ell}^{(i,j)}\|_{\mathbb{F}}^2 \leq O(2^{2k/L} dmL + DL). \quad (8)$$

340 Making use of the 1-homogeneous property of the ReLU function, by scaling all the weights in
 341 the neural network, these building blocks can be combined into a neural network with residual
 342 connections, that approximate the target function and satisfy our constraint on the norm of weights.
 343 See [Section C.4](#) for the detail.

344 By applying [Lemma 12](#), which shows that any L -layer feedforward neural network can be reformu-
 345 lated as an $L + L_0 - 1$ -layer convolution neural network, the neural network constructed above can
 346 be converted into a ConvResNeXt that satisfies the conditions in [Theorem 1](#).

347 B.2 Estimation error

348 The formal theorem for the upper bound of estimation error of \hat{f} is presented as follows:

Theorem 5. Suppose Assumption 1 holds. Set $L = L' + L_0 - 1$, $L' \geq 3$, where $L_0 = \lceil \frac{D}{K-1} \rceil$, and

$$MN \geq C_{\mathcal{M}}P, \quad P = O(n^{\frac{1-2/L}{2\alpha/d(1-1/L)+1-2/pL}}), \quad w \geq C_1(dm + D).$$

349 Let \hat{f} be the global minimizer given in (1) with the function class $\mathcal{F} =$
 350 $\mathcal{F}^{\text{Conv}}(N, M, L, K, w, B_{\text{res}}, B_{\text{out}})$. Then we have

$$\begin{aligned} \mathbb{E}_{\mathcal{D}}[\text{Loss}(\hat{f}(x), y)] &\leq \mathbb{E}_{\mathcal{D}}[\text{Loss}(f^*(x), y)] + C_7 \left(\frac{K^{-\frac{2}{L-2}} w^{\frac{3L-4}{L-2}} L^{\frac{3L-2}{L-2}}}{n} \right)^{\frac{\alpha/d(1-2/L)}{2\alpha/d(1-1/L)+1-2/(pL)}} \\ &\quad + C_8 \exp(-C_6 L'), \end{aligned}$$

351 where the logarithmic terms are omitted. C_1 is the constant defined in Theorem 1, C_7, C_8 are
 352 constants that depend on $C_{\mathbb{F}}, C_{\mathcal{M}}, d, m, K$ is the size of the convolution kernel.

353 To prove the above theorem, we first compute the covering number of an overparameterized ConvResNeXt
 354 with norm-constraint as in Lemma 6, then compute the critical radius of this function class
 355 using the covering number as in Corollary 19. The critical radius can be used to bound the estimation
 356 error as in Theorem 14.20 in Wainwright [23]. The proof is deferred to Section D.2.

357 **Lemma 6.** Consider a neural network defined in Definition 1. Let the last layer of this neural network
 358 is a single linear layer with norm $\|W_{\text{out}}\|_{\mathbb{F}}^2 \leq B_{\text{out}}$. Let the input of this neural network satisfy
 359 $\|\mathbf{x}\|_2 \leq 1, \forall x$, and is concatenated with 1 before feeding into this neural network so that part of the
 360 weight plays the role of the bias. The covering number of this neural network is bounded by

$$\log \mathcal{N}(\cdot, \delta) \lesssim w^2 L B_{\text{res}}^{\frac{1}{1-2/L}} K^{\frac{2-2/L}{1-2/L}} (B_{\text{out}}^{1/2} \exp((KB_{\text{res}}/L)^{L/2}))^{\frac{2/L}{1-2/L}} \delta^{-\frac{2/L}{1-2/L}}, \quad (9)$$

361 where the logarithmic term is omitted.

362 The key idea of the proof is to split the building block into two types (“small blocks” and “large
 363 blocks”) depending on whether the total norm of the weights in the building block is smaller than
 364 ϵ or not. By properly choosing ϵ , we prove that if all the “small blocks” in this neural network are
 365 removed, the perturbation to the output for any input $\|\mathbf{x}\| \leq 1$ is no more than $\delta/2$, so the covering
 366 number of the ConvResNeXt is only determined by the number of “large blocks”, which is no more
 367 than B_{res}/ϵ .

Proof. Using the inequality of arithmetic and geometric means, from Proposition 20, Proposition 22
 and Proposition 23, if any residual block is removed, the perturbation to the output is no more than

$$(KB_m/L)^{L/2} B_{\text{out}}^{1/2} \exp((KB_{\text{res}}/L)^{L/2}),$$

368 where B_m is the total norm of parameters in this block. Because of that, the residual blocks can be
 369 divided into two kinds depending on the norm of the weights $B_m < \epsilon$ (“small blocks”) and $B_m \geq \epsilon$
 370 (“large blocks”). If all the “small blocks” are removed, the perturbation to the output for any input
 371 $\|\mathbf{x}\|_2 \leq 1$ is no more than

$$\begin{aligned} &\exp((KB_{\text{res}}/L)^{L/2}) B_{\text{out}}^{1/2} \sum_{m: B_m < \epsilon} (KB_m/L)^{L/2} \\ &\leq \exp((KB_{\text{res}}/L)^{L/2}) B_{\text{out}}^{1/2} \sum_{m: B_m < \epsilon} (KB_m/L)(K\epsilon/L)^{L/2-1} \\ &\leq \exp((KB_{\text{res}}/L)^{L/2}) K^{L/2} B_{\text{res}} B_{\text{out}}^{1/2} (\epsilon/L)^{L/2-1}/L. \end{aligned}$$

372 Choosing $\epsilon = L \left(\frac{\delta L}{2 \exp((B_{\text{res}}/L)^{L/2}) K^{L/2} B_{\text{res}} B_{\text{out}}^{1/2}} \right)^{\frac{1}{L/2-1}}$, the perturbation above is no more than
 373 $\delta/2$. The covering number can be determined by the number of the “large blocks” in the neural
 374 network, which is no more than B_{res}/ϵ .

375 Taking our choice of ϵ into Proposition 13 and noting that for any block, $B_{\text{in}} L_{\text{post}} \leq$
 376 $B_{\text{out}}^{1/2} \exp((KB_{\text{res}}/L)^{L/2})$ finishes the proof, where B_{in} is the upper bound of the input to this
 377 block as defines in Proposition 13, and L_{post} is the Lipschitz parameter of all the layers following
 378 the block.

379 \square

380 **Remark 1.** *The proof of Lemma 6 shows that under weight decay, the building blocks in a ConvResNeXt are sparse, i.e. only a finite number of blocks contribute non-trivially to the network even*
 381 *though the model can be overparameterized. This explains why a ConvResNeXt can generalize well*
 382 *despite overparameterization, and provide a new perspective in explaining why residual connections*
 383 *improve the performance of deep neural networks.*

385 C Proof of the approximation theory

386 C.1 Decompose the target function into the sum of locally supported functions.

387 **Lemma 7.** *Approximating Besov function on a smooth manifold using B-spline: Let $f \in B_{p,q}^\alpha(\mathcal{M})$.*
 388 *There exists a decomposition of f :*

$$f(\mathbf{x}) = \sum_{i=1}^{C_{\mathcal{M}}} \tilde{f}_i \circ \phi_i(\mathbf{x}) \times \mathbf{1}(\mathbf{x} \in B(\mathbf{c}_i, r)),$$

389 *and $\tilde{f}_i = f \cdot \rho_i \in B_{p,q}^\alpha$, $\sum_{i=1}^{C_{\mathcal{M}}} \|\tilde{f}_i\|_{B_{p,q}^\alpha} \leq C\|f\|_{B_{p,q}^\alpha(\mathcal{M})}$, $\phi_i : \mathcal{M} \rightarrow \mathbb{R}^d$ are linear projections,*
 390 *$B(\mathbf{c}_i, r)$ denotes the unit ball with radius r and center \mathbf{c}_i .*

391 The lemma is inferred by the existence of the partition of unity, which is given in Proposition 3.

392 C.2 Locally approximate the decomposed functions using cardinal B-spline basis functions.

393 **Proposition 8.** *For any function in the Besov space on a compact smooth manifold $f^* \in B_{p,q}^s(\mathcal{M})$,*
 394 *any $N \geq 0$, there exists an approximated to f^* using cardinal B-spline basis functions:*

$$\tilde{f} = \sum_{i=1}^{C_{\mathcal{M}}} \sum_{j=1}^P a_{i,k_j,s_j} M_{m,k_j,s_j} \circ \phi_i \times \mathbf{1}(\mathbf{x} \in B(\mathbf{c}_i, r)),$$

395 *where m is the integer satisfying $0 < \alpha < \min(m, m - 1 + 1/p)$, $M_{m,k,s} = M_m(2^k(\cdot - s))$, M_m*
 396 *denotes the B-spline basis function defined in (4), the approximation error is bounded by*

$$\|f - \tilde{f}\|_\infty \leq C_9 C_{\mathcal{M}} P^{-\alpha/d}$$

397 *and the coefficients satisfy*

$$\|\{2^{kj} a_{i,k_j,s_j}\}_{i,j}\|_p \leq C_{10} \|f\|_{B_{p,q}^\alpha(\mathcal{M})}$$

398 *for some constant C_9, C_{10} that only depends on α .*

399 As will be shown below, the scaled coefficients $2^{kj} a_{i,k_j,s_j}$ corresponds to the total norm of the
 400 parameters in the neural network to approximate the B-spline basis function, so this lemma is the key
 401 to get the bound of norm of parameters in (10).

402 *Proof.* From the definition of $B_{p,q}^\alpha(\mathcal{M})$, and applying Proposition 3, there exists a decomposition of
 403 f^* as

$$f^* = \sum_{i=1}^{C_{\mathcal{M}}} (f_i) = \sum_{i=1}^{C_{\mathcal{M}}} (f_i \circ \phi_i^{-1}) \circ \phi_i \times \mathbf{1}_{U_i},$$

404 where $f_i := f^* \cdot \rho_i$, ρ_i satisfy the condition in Definition 6, and $f_i \circ \phi_i^{-1} \in B_{p,q}^\alpha$. Using Proposition 16,
 405 for any i , one can approximate $f_i \circ \phi_i^{-1}$ with \bar{f}_i :

$$\bar{f}_i = \sum_{j=1}^P a_{i,k_j,s_j} M_{m,k_j,s_j}$$

such that $\|f_i \circ \phi_i^{-1}\|_\infty \leq C_1 M^{-\alpha/d}$, and the coefficients satisfy

$$\|\{2^{kj} a_{i,k_j,s_j}\}_j\|_p \leq C_{10} \|f_i \circ \phi_i^{-1}\|_{B_{p,q}^\alpha}.$$

Define

$$\bar{f} = \sum_{i=1}^{C_{\mathcal{M}}} \bar{f}_i \circ \phi_i \times \mathbf{1}_{U_i}.$$

406 one can verify that $\|f - \tilde{f}\|_\infty \leq C_9 C_{\mathcal{M}} N^{-\alpha/d}$. On the other hand, using triangular inequality (and
 407 padding the vectors with 0),

$$\|\{2^{kj} a_{i,k_j,s_j}\}_{i,j}\|_p \leq \sum_{i=1}^{C_{\mathcal{M}}} \|\{2^{kj} a_{i,k_j,s_j}\}_j\|_p \leq \sum_{i=1}^{C_{\mathcal{M}}} C_{10} \|f_i \circ \phi_i^{-1}\|_{B_{p,q}^\alpha} = C_{10} \|f^*\|_{B_{p,q}^\alpha(\mathcal{M})},$$

408 which finishes the proof.

409

□

410 C.3 Neural network for chart selection

411 In this section, we demonstrate that a feedforward neural network can approximate the chart selection
 412 function $z \times \mathbf{1}(\mathbf{x} \in B(\mathbf{c}_i, r))$, and it is error-free as long as $z = 0$ when $r < d(\mathbf{x}, \mathbf{c}_i) < R$. We start
 413 by proving the following supporting lemma:

414 **Proposition 9.** *Fix some constant $B > 0$. For any $\mathbf{x}, \mathbf{c} \in \mathbb{R}^D$ satisfying $|x_i| \leq B$ and $|c_i| \leq B$ for
 415 $i = 1, \dots, D$, there exists an L -layer neural network $\tilde{d}(\mathbf{x}; \mathbf{c})$ with width $w = O(d)$ that approximates
 416 $d^2(\mathbf{x}; \mathbf{c}) = \sum_{i=1}^D (x_i - c_i)^2$ such that $|\tilde{d}^2(\mathbf{x}; \mathbf{c}) - d^2(\mathbf{x}; \mathbf{c})| \leq 8DB^2 \exp(-C_{11}L)$ with an absolute
 417 constant $C_{11} > 0$ when $d(\mathbf{x}; \mathbf{c}) < \tau$, and $\tilde{d}^2(\mathbf{x}; \mathbf{c}) \geq \tau^2$ when $d(\mathbf{x}; \mathbf{c}) \geq \tau$, and the norm of the
 418 neural network is bounded by*

$$\sum_{\ell=1}^L \|W_\ell\|_F^2 + \|b_\ell\|_2^2 \leq C_{12}DL.$$

419 *Proof.* The proof is given by construction. By Proposition 2 in Yarotsky(2017), the function $f(x) =$
 420 x^2 on the segment $[0, 2B]$ can be approximated with any error $\epsilon > 0$ by a ReLU network g having
 421 depth and the number of neurons and weight parameters no more than $c \log(4B^2/\epsilon)$ with an absolute
 422 constant c . The width of the network g is an absolute constant. We also consider a single layer ReLU
 423 neural network $h(t) = \text{ReLU}(t) - \text{ReLU}(-t)$, which is equal to the absolute value of the input.

424 Now we consider a neural network $G(\mathbf{x}; \mathbf{c}) = \sum_{i=1}^D g \circ h(x_i - c_i)$. Then for any $\mathbf{x}, \mathbf{c} \in \mathbb{R}^D$
 425 satisfying $|x_i| \leq B$ and $|c_i| \leq B$ for $i = 1, \dots, D$, we have

$$\begin{aligned} |G(\mathbf{x}; \mathbf{c}) - d^2(\mathbf{x}; \mathbf{c})| &\leq \left| \sum_{i=1}^D g \circ h(x_i - c_i) - \sum_{i=1}^D (x_i - c_i)^2 \right| \\ &\leq \sum_{i=1}^D |g \circ h(x_i - c_i) - (x_i - c_i)^2| \\ &\leq D\epsilon. \end{aligned}$$

426 Moreover, define another neural network

$$\begin{aligned} F(\mathbf{x}; \mathbf{c}) &= -\text{ReLU}(\tau^2 - D\epsilon - G(\mathbf{x}; \mathbf{c})) + \tau^2 \\ &= \begin{cases} G(\mathbf{x}; \mathbf{c}) + D\epsilon & \text{if } G(\mathbf{x}; \mathbf{c}) < \tau^2 - D\epsilon, \\ \tau^2 & \text{if } G(\mathbf{x}; \mathbf{c}) \geq \tau^2 - D\epsilon, \end{cases} \end{aligned}$$

427 which has depth and the number of neurons no more than $c' \log(4B^2/\epsilon)$ with an absolute constant c' .
 428 The weight parameters of G are upper bounded by $\max\{\tau^2, D\epsilon, c \log(4B^2/\epsilon)\}$ and the width of G
 429 is $O(D)$.

430 If $d^2(\mathbf{x}; \mathbf{c}) < \tau^2$, we have

$$\begin{aligned} |F(\mathbf{x}; \mathbf{c}) - d^2(\mathbf{x}; \mathbf{c})| &= |-\text{ReLU}(\tau^2 - D\epsilon - G(\mathbf{x}; \mathbf{c})) + \tau^2 - d^2(\mathbf{x}; \mathbf{c})| \\ &= \begin{cases} |G(\mathbf{x}; \mathbf{c}) - d^2(\mathbf{x}; \mathbf{c}) + D\epsilon| & \text{if } G(\mathbf{x}; \mathbf{c}) < \tau^2 - D\epsilon, \\ \tau^2 - d^2(\mathbf{x}; \mathbf{c}) & \text{if } G(\mathbf{x}; \mathbf{c}) \geq \tau^2 - D\epsilon. \end{cases} \end{aligned}$$

431 For the first case when $G(\mathbf{x}; \mathbf{c}) < \tau^2 - D\epsilon$, $|F(\mathbf{x}; \mathbf{c}) - d^2(\mathbf{x}; \mathbf{c})| \leq 2D\epsilon$ since $d^2(\mathbf{x}; \mathbf{c})$ can be
 432 approximated by $G(\mathbf{x}; \mathbf{c})$ up to an error ϵ . For the second case when $G(\mathbf{x}; \mathbf{c}) \geq \tau^2 - D\epsilon$, we have
 433 $d^2(\mathbf{x}; \mathbf{c}) \geq G(\mathbf{x}; \mathbf{c}) - D\epsilon \geq \tau^2 - 2D\epsilon$ and . Thereby we also have $|F(\mathbf{x}; \mathbf{c}) - d^2(\mathbf{x}; \mathbf{c})| \leq 2D\epsilon$.

434 If $d^2(\mathbf{x}; \mathbf{c}) \geq \tau^2$ instead, we will obtain $G(\mathbf{x}; \mathbf{c}) \geq d^2(\mathbf{x}; \mathbf{c}) - D\epsilon \geq \tau^2 - D\epsilon$. This gives that
 435 $F(\mathbf{x}; \mathbf{c}) = \tau^2$ in this case.

436 Finally, we take $\epsilon = 4B^2 \exp(-L/c')$. Then $F(\mathbf{x}; \mathbf{c})$ is an L -layer neural network with $O(L)$
 437 neurons. The weight parameters of G are upper bounded by $\max\{\tau^2, 4DB^2 \exp(-L/c'), cL/c'\}$
 438 and the width of G is $O(D)$. Moreover, $F(\mathbf{x}; \mathbf{c})$ satisfies $|F(\mathbf{x}; \mathbf{c}) - d^2(\mathbf{x}; \mathbf{c})| < 8DB^2 \exp(-L/c')$
 439 if $d^2(\mathbf{x}; \mathbf{c}) \leq \tau^2$ and $F(\mathbf{x}; \mathbf{c}) = \tau^2$ if $d^2(\mathbf{x}; \mathbf{c}) \geq \tau^2$. \square

440 **Proposition 10.** *There exists a single layer ReLU neural network that approximates \tilde{x} , such that for*
 441 *all $0 \leq x \leq C, y \in \{0, 1\}$, $x \tilde{x} y = x$ when $y = 1$, and $x \tilde{x} y = 0$ when either $x = 0$ or $y = 0$.*

442 *Proof.* Consider a single layer neural network $g(\mathbf{x}, y) := A_2 \text{ReLU}(A_1(\mathbf{x}, y)^\top)$ with no bias, where

$$A_1 = \begin{bmatrix} -\frac{1}{C} & 1 \\ 0 & 1 \end{bmatrix}, \quad A_2 = \begin{bmatrix} -C \\ C \end{bmatrix}.$$

443 Then we can rewrite the neural network g as $g(x, y) = -C \text{ReLU}(-x/C + y) + C \text{ReLU}(y)$. If
 444 $y = 1$, we will have $g(x, y) = -C \text{ReLU}(-x/C + 1) + C = x$, since $x \leq C$. If $y = 0$, we will
 445 have $g(x, y) = -C \text{ReLU}(-x/C) = 0$, since $x \geq 0$. Thereby we can conclude the proof. \square

446 By adding a single linear layer

$$y = \frac{1}{R - r - 2\Delta} (\text{ReLU}(R - \Delta - x) - \text{ReLU}(r + \Delta - x))$$

447 after the one shown in [Proposition 9](#), where $\Delta = 8DB^2 \exp(-CL)$ denotes the error in [Proposition 9](#),
 448 one can approximate the indicator function $\mathbf{1}(x \in B(\mathbf{c}_i, r))$ such that it is error-free when $d(\mathbf{x}, \mathbf{c}_i) \leq$
 449 r or $\geq R$. Choosing $R \leq \tau/2, r < R - 2\Delta$, and combining with [Proposition 10](#), the proof is finished.
 450 Considering that f_i is locally supported on $B(\mathbf{c}_i, r)$ for all i by our method of construction, the chart
 451 selection part does not incur any error in the output.

452 C.4 Constructing the neural network to approximate the target function

453 In this section, we focus on the neural network with the same architecture as a ResNeXt in [Definition 1](#)
 454 but replacing each building block with a feedforward neural network, and prove that it can achieve
 455 the same approximation error as in [Theorem 1](#). For technical simplicity, we assume that the target
 456 function $f^* \in [0, 1]$ without loss of generality. Then our analysis automatically holds for any bounded
 457 function.

458 **Theorem 11.** *For any f^* under the same condition as [Theorem 1](#), any neural network architecture*
 459 *with residual connections containing N number of residual blocks and each residual block contains*
 460 *M number of feedforward neural networks in parallel, where the depth of each feedforward neural*
 461 *networks is L , width is w :*

$$f = \mathbf{W}_{\text{out}} \cdot \left(1 + \sum_{m=1}^M f_{N,m} \right) \circ \cdots \circ \left(1 + \sum_{m=1}^M f_{1,m} \right)$$

$$f_{n,m} = \mathbf{W}_L^{(n,m)} \text{ReLU}(\mathbf{W}_{L-1}^{(n,m)} \dots \text{ReLU}(\mathbf{W}_1^{(n,m)} \mathbf{x})) \circ P(\mathbf{x}),$$

462 where $P(\mathbf{x}) = [\mathbf{x}^T, 1, 0]^T$ is the padding operation,

463 satisfying

$$MN \geq C_{\mathcal{M}} P, \quad w \geq C_1(dm + D),$$

$$B_{\text{res}} := \sum_{n=1}^N \sum_{m=1}^M \sum_{\ell=1}^L \|\mathbf{W}_\ell^{(n,m)}\|_{\text{F}}^2 \leq C_2 L, \quad (10)$$

$$B_{\text{out}} := \|\mathbf{W}_{\text{out}}\|_{\text{F}}^2 \leq C_3 C_{\text{F}}^2 ((dm + D)L)^L (C_{\mathcal{M}} P)^{L-2/p},$$

464 there exists an instance f of this ResNeXt class, such that

$$\|f - f^*\|_{\infty} \leq C_{\text{F}} C_{\mathcal{M}} \left(C_4 P^{-\alpha/d} + C_5 \exp(-C_6 L \log P) \right), \quad (11)$$

465 where $C_1, C_2, C_3, C_4, C_5, C_6$ are the same constants as in [Theorem 1](#).

466 *Proof.* We first construct a parallel neural network to approximate the target function, then scale the
 467 weights to meet the norm constraint while keeping the model equivalent to the one constructed in the
 468 first step, and finally transform this parallel neural network into the ConvResNeXt as claimed.

469 Combining Lemma 17, Proposition 9 and Proposition 10, by putting the neural network in Lemma 17
 470 and Proposition 9 in parallel and adding the one in Proposition 10 after them, one can construct
 471 a feedforward neural network with bias with depth L , width $w = O(d) + O(D) = O(d)$, that
 472 approximates $M_{m,k_j,s_j}(\mathbf{x}) \times \mathbf{1}(\mathbf{x} \in B(\mathbf{c}_i, r))$ for any i, j .

473 To construct the neural network with residual connections that approximates f^* , we follow the
 474 method in Oono and Suzuki [16], Liu et al. [14]. This network uses separate channels for the inputs
 475 and outputs. Let the input to one residual layer be $[\mathbf{x}_1, y_1]$, the output is $[\mathbf{x}_1, y_1 + f(\mathbf{x}_1)]$. As a result,
 476 if one scale the outputs of all the building blocks by any scalar a , then the last channel of the output
 477 of the entire network is also scaled by a . This property allows us to scale the weights in each building
 478 block while keeping the model equivalent. To compensate for the bias term, Proposition 18 can be
 479 applied. This only increases the total norm of each building block by no larger than a constant term
 480 that depends only L , which is no more than a factor of constant.

481 Let the neural network constructed above has parameter $\tilde{\mathbf{W}}_1^{(i,j)}, \tilde{\mathbf{b}}_1^{(i,j)}, \dots, \tilde{\mathbf{W}}_L^{(i,j)}, \mathbf{b}_L^{(i,j)}$ in each
 482 layer, one can construct a building block without bias as

$$\mathbf{W}_1^{(i,j)} = \begin{bmatrix} \tilde{\mathbf{W}}_1^{(i,j)} & \tilde{\mathbf{b}}_1^{(i,j)} & 0 \\ 0 & 1 & 0 \end{bmatrix}, \quad \mathbf{W}_\ell^{(i,j)} = \begin{bmatrix} \tilde{\mathbf{W}}_\ell^{(i,j)} & \tilde{\mathbf{b}}_\ell^{(i,j)} \\ 0 & 1 \end{bmatrix}, \quad \mathbf{W}_L^{(i,j)} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ \tilde{\mathbf{W}}_L^{(i,j)} & \tilde{\mathbf{b}}_L^{(i,j)} \end{bmatrix}.$$

483 Remind that the input is padded with the scalar 1 before feeding into the neural network, the above
 484 construction provide an equivalent representation to the neural network including the bias, and route
 485 the output to the last channel. From Lemma 17, it can be seen that the total square norm of this block
 486 is bounded by (8).

487 Finally, we scale the weights in the each block, including the “1” terms to meet the norm constraint.
 488 Thanks to the 1-homogeneous property of ReLU layer, and considering that the network we construct
 489 use separate channels for the inputs and outputs, the model is equivalent after scaling. Actually the
 490 property above allows the tradeoff between \bar{B}_{res} and B_{out} . If all the weights in the residual blocks are
 491 scaled by an arbitrary positive constant c , and the weight in the last layer \mathbf{W}_{out} is scaled by c^{-L} , the
 492 model is still equivalent. We only need to scale the all the weights in this block with $|a_{i,k_j,s_j}|^{1/L}$,
 493 setting the sign of the weight in the last layer as $\text{sign}(a_{i,k_j,s_j})$, and place $C_{\mathcal{M}}P$ number of these
 494 building blocks in this neural network with residual connections. Since this block always output 0
 495 in the first $D + 1$ channels, the order and the placement of the building blocks does not change the
 496 output. The last fully connected layer can be simply set to

$$\mathbf{W}_{\text{out}} = [0, \dots, 0, 1], b_{\text{out}} = 0.$$

497 Combining Proposition 16 and Lemma 15, the norm of this ResNeXt we construct satisfy

$$\begin{aligned} \bar{B}_{\text{res}} &\leq \sum_{i=1}^{C_{\mathcal{M}}} \sum_{j=1}^P a_{i,k_j,s_j}^{2/L} (2^{2k/L} C_{14} dmL + C_{12} DL) \\ &\leq \sum_{i=1}^{C_{\mathcal{M}}} \sum_{j=1}^P (2^k a_{i,k_j,s_j})^{2/L} (C_{14} dmL + C_{12} DL) \\ &\leq (C_{\mathcal{M}}P)^{1-2/(pL)} \|\{2^k a_{i,k_j,s_j}\}\|_p^{2/L} (C_{14} dmL + C_{12} DL) \\ &\leq (C_{10} C_{\mathbb{F}})^{2/L} (C_{\mathcal{M}}P)^{1-2/(pL)} (C_{14} dmL + C_{12} DL), \\ \bar{B}_{\text{out}} &\leq 1. \end{aligned}$$

498 By scaling all the weights in the residual blocks by $\bar{B}_{\text{res}}^{-1/2}$, and scaling the output layer by $\bar{B}_{\text{res}}^{L/2}$, the
 499 network that satisfy (10) can be constructed. \square

500 Notice that the chart selection part does not introduce error by our way of construction, we only
 501 need to sum over the error in Section B.1 and Section B.1, and notice that for any \mathbf{x} , for any linear
 502 projection ϕ_i , the number of B-spline basis functions $M_{m,k,s}$ that is nonzero on \mathbf{x} is no more than
 503 $m^d \log P$, the approximation error of the constructed neural network can be proved.

504 **C.5 Constructing a convolution neural network to approximate the target function**

505 In this section, we prove that any feedforward neural network can be realized by a convolution neural
 506 network with similar size and norm of parameters. The proof is similar to Theorem 5 in [16].

507 **Lemma 12.** *For any feedforward neural network with depth L' , width w' , input dimension h and*
 508 *output dimension h' , for any kernel size $K > 1$, there exists a convolution neural network with depth*
 509 *$L = L' + L_0 - 1$, where $L_0 = \lceil \frac{h-1}{K-1} \rceil$ number of channels $w = 4w'$, and the first dimension of*
 510 *the output equals the output of the feedforward neural network for all inputs, and the norm of the*
 511 *convolution neural network is bounded as*

$$\sum_{\ell=1}^L \|\mathbf{W}_\ell\|_{\mathbb{F}}^2 \leq 4 \sum_{\ell=1}^{L'} \|\mathbf{W}'_\ell\|_{\mathbb{F}}^2 + 4w'L_0,$$

512 where $\mathbf{W}'_1 \in \mathbb{R}^{w' \times h'}$; $\mathbf{W}'_\ell \in \mathbb{R}^{w' \times w'}$, $\ell = 2, \dots, L' - 1$; $\mathbf{W}'_{L'} \in \mathbb{R}^{h' \times w'}$ are the weights in the
 513 feedforward neural network, and $\mathbf{W}_1 \in \mathbb{R}^{K \times w \times h}$, $\mathbf{W}_\ell \in \mathbb{R}^{K \times w \times w}$, $\ell = 2, \dots, L - 1$; $\mathbf{W}_L \in$
 514 $\mathbb{R}^{K \times h \times w}$ are the weights in the convolution neural network.

515 *Proof.* We follow the same method as Oono and Suzuki [16] to construct the CNN that is equivalent
 516 to the feedforward neural network. By combining Oono and Suzuki [16] lemma 1 and lemma 2, for
 517 any linear transformation, one can construct a convolution neural network with at most $L_0 = \lceil \frac{h-1}{K-1} \rceil$
 518 convolution layers and 4 channels, where h is the dimension of input, which equals $D + 1$ in our
 519 case, such that the first dimension in the output equals the linear transformation, and the norm of all
 520 the weights is no more than

$$\sum_{\ell=1}^{L_0} \|\mathbf{W}_\ell\|_{\mathbb{F}}^2 \leq 4L_0, \quad (12)$$

521 where \mathbf{W}_ℓ is the weight of the linear transformation. Putting w number of such convolution neural
 522 networks in parallel, a convolution neural network with L_0 layers and $4w$ channels can be constructed
 523 to implement the first layer in the feedforward neural network.

524 To implement the remaining layers, one choose the convolution kernel $\mathbf{W}_{\ell+L_0-1}[:, i, j] =$
 525 $[0, \dots, \mathbf{W}'_\ell[i, j], \dots, 0], \forall 1 \leq i, j \leq w$, and pad the remaining parts with 0, such that this con-
 526 volution layer is equivalent to the linear layer applied on the dimension of channels. Noticing that
 527 this conversion does not change the norm of the parameters in each layer. Adding both sides of (12)
 528 by the norm of the $2 - L'$ -th layer in both models finishes the proof. \square

529 **D Proof of the estimation theory**

530 **D.1 Covering number of a neural network block**

531 **Proposition 13.** *If the input to a ReLU neural network is bounded by $\|\mathbf{x}\|_2 \leq B_{\text{in}}$, the covering*
 532 *number of the ReLU neural network defined in Proposition 20 is bounded by*

$$\mathcal{N}(\mathcal{F}_{NN}, \delta, \|\cdot\|_2) \leq \left(\frac{B_{\text{in}}(B/L)^{L/2} w L}{\delta} \right)^{w^2 L}.$$

533 *Proof.* Similar to Proposition 20, we only consider the case $\|W_\ell\|_{\mathbb{F}} \leq \sqrt{B/L}$. For any $1 \leq \ell \leq L$,
 534 for any $W_1, \dots, W_{\ell-1}, W_\ell, W'_\ell, W_{\ell+1}, \dots, W_L$ that satisfy the above constraint and $\|W_\ell - W'_\ell\|_{\mathbb{F}} \leq \epsilon$,
 535 define $g(\dots; W_1, \dots, W_L)$ as the neural network with parameters W_1, \dots, W_L , we can see

$$\begin{aligned} & \|g(\mathbf{x}; W_1, \dots, W_{\ell-1}, W_\ell, W_{\ell+1}, \dots, W_L) - g(\mathbf{x}; W_1, \dots, W_{\ell-1}, W_\ell, W_{\ell+1}, \dots, W_L)\|_2 \\ & \leq (B/L)^{(L-\ell)/2} \|W_\ell - W'_\ell\|_2 \|ReLU(W_{\ell-1} \dots ReLU(W_1(\mathbf{x})))\|_2 \\ & \leq (B/L)^{(L-1)/2} B_{\text{in}} \epsilon. \end{aligned}$$

536 Choosing $\epsilon = \frac{\delta}{L(B/L)^{(L-1)/2}}$, the above inequality is no larger than δ/L . Taking the sum over ℓ , we
 537 can see that for any $W_1, W'_1, \dots, W_L, W'_L$ such that $\|W_\ell - W'_\ell\|_{\mathbb{F}} \leq \epsilon$,

$$\|g(\mathbf{x}; W_1, \dots, W_L) - g(\mathbf{x}; W'_1, \dots, W'_L)\|_2 \leq \delta.$$

538 Finally, observe that the covering number of W_ℓ is bounded by

$$\mathcal{N}(\{W : \|W\|_F \leq B\}, \epsilon, \|\cdot\|_F) \leq \left(\frac{2Bw}{\epsilon}\right)^{w^2}. \quad (13)$$

539 Substituting B and ϵ and taking the product over ℓ finishes the proof. \square

540 **Proposition 14.** *If the input to a ReLU convolution neural network is bounded by $\|x\|_2 \leq B_{\text{in}}$, the*
 541 *covering number of the ReLU neural network defined in Definition 1 is bounded by*

$$\mathcal{N}(\mathcal{F}_{\text{NN}}, \delta, \|\cdot\|_2) \leq \left(\frac{B_{\text{in}}(BK/L)^{L/2}wL}{\delta}\right)^{w^2KL}.$$

542 *Proof.* Similar to Proposition 13, for any $1 \leq \ell \leq L$, for any $W_1, \dots, W_{\ell-1}, W_\ell, W'_\ell, W_{\ell+1}, \dots, W_L$
 543 that satisfy the above constraint and $\|W_\ell - W'_\ell\|_F \leq \epsilon$, define $g(\dots; W_1, \dots, W_L)$ as the neural
 544 network with parameters W_1, \dots, W_L , we can see

$$\begin{aligned} & \|g(\mathbf{x}; W_1, \dots, W_{\ell-1}, W_\ell, W_{\ell+1}, \dots, W_L) - g(\mathbf{x}; W_1, \dots, W_{\ell-1}, W_\ell, W_{\ell+1}, \dots, W_L)\|_2 \\ & \leq K^{L/2}(B/L)^{(L-\ell)/2} \|W_\ell - W'_\ell\|_2 \|ReLU(W_{\ell-1} \dots ReLU(W_1(\mathbf{x})))\|_2 \\ & \leq K^{L/2}(B/L)^{(L-1)/2} B_{\text{in}} \epsilon, \end{aligned}$$

545 where the first inequality comes from Proposition 24. Choosing $\epsilon = \frac{\delta}{K^{L/2}B_{\text{in}}L(B/L)^{(L-1)/2}}$, the
 546 above inequality is no larger than δ/L . Taking this into (13) finishes the proof. \square

547 D.2 Proof of Theorem 5

548 Define $\tilde{f} = \arg \min_f \mathbb{E}_{\mathcal{D}}[\text{Loss}(f)]$. From Theorem 14.20 in Wainwright [23], for any function class
 549 $\partial\mathcal{F}$ that is star-shaped around \tilde{f} , the empirical risk minimizer $\hat{f} = \arg \min_{f \in \mathcal{F}} \text{Loss}_n(f)$ satisfy

$$\mathbb{E}_{\mathcal{D}}[\text{Loss}(\hat{f})] \leq \mathbb{E}_{\mathcal{D}}[\text{Loss}(\tilde{f})] + 10\delta_n(2 + \delta_n) \quad (14)$$

550 with probability at least $1 - c_1 \exp(-c_2 n \delta_n^2)$ for any δ_n that satisfy (18), where c_1, c_2 are universal
 551 constants.

552 The function of neural networks is not star-shaped, but can be covered by a star-shaped function class.
 553 Specifically, let $\{f - \tilde{f} : f \in \mathcal{F}^{\text{Conv}}\} \subset \{f_1 - f_2 : f_1, f_2 \in \mathcal{F}^{\text{Conv}}\} := \partial\mathcal{F}$.

554 Any function in $\partial\mathcal{F}$ can be represented using a ResNeXt: one can put two neural networks of the
 555 same structure in parallel, adjusting the sign of parameters in one of the neural networks and summing
 556 up the result, which increases M, B_{res} and B_{out} by a factor of 2. This only increases the log covering
 557 number in (9) by a factor of constant (remind that $B_{\text{res}} = O(1)$ by assumption).

558 Taking the log covering number of the ResNeXt (9), the sufficient condition for the critical radius as
 559 in (18) is

$$\begin{aligned} n^{-1/2}wL^{1/2}B_{\text{res}}^{\frac{1}{2-4/L}}K^{\frac{1-1/L}{1-2/L}}(B_{\text{out}}^{1/2}\exp((KB_{\text{res}}/L)^{L/2}))^{\frac{1/L}{1-2/L}}\delta_n^{\frac{1-3/L}{1-2/L}} & \lesssim \frac{\delta_n^2}{4}, \\ \delta_n \gtrsim K(w^2L)^{\frac{1-2/L}{2-2/L}}B_{\text{res}}^{\frac{1}{2-2/L}}(B_{\text{out}}^{1/2}\exp((KB_{\text{res}}/L)^{L/2}))^{\frac{1/L}{1-1/L}}n^{-\frac{1-2/L}{2-2/L}}, \end{aligned} \quad (15)$$

560 where \lesssim hides the logarithmic term.

561 Because Loss is 1-Lipschitz, we have

$$\text{Loss}(f) \leq \text{Loss}(\tilde{f}) + \|f - \tilde{f}\|_\infty.$$

Choosing

$$P = O\left(\left(\frac{K^{-\frac{2}{L-2}}w^{\frac{3L-4}{L-2}}L^{\frac{3L-2}{L-2}}}{n}\right)^{-\frac{1-2/L}{2\alpha/d(1-1/L)+1-2/pL}}\right),$$

562 and taking Theorem 1 and (15) into (14) finishes the proof.

563 E Lower bound of error

564 In this section, we study the minimax lower bound of any estimator for Besov functions on a d -
 565 dimensional manifold. It suffices to consider the manifold \mathcal{M} as a d -dimensional hypersurface.
 566 Without the loss of generalization, assume that $\frac{\partial \text{Loss}(y)}{\partial y} \geq 0.5$ for $-\epsilon \leq y \leq \epsilon$. Define the function
 567 space

$$\mathcal{F} = \left\{ f = \sum_{j_1, \dots, j_d=1}^s \pm \frac{\epsilon}{s^\alpha} \times M^{(m)}((\mathbf{x} - \mathbf{j})/s) \right\}, \quad (16)$$

568 where $M^{(m)}$ denotes the Cardinal B-spline basis function that is supported on $(0, 1)^d$, $\mathbf{j} =$
 569 $[j_1, \dots, j_d]$. The support of each B-spline basis function splits the space into s^d number of blocks,
 570 where the target function in each block has two choices (positive or negative), so the total number of
 571 different functions in this function class is $|\mathcal{F}| = 2^{s^d}$. Using Dũng [5, Theorem 2.2], we can see that
 572 for any $f \in \mathcal{F}$,

$$\|f\|_{B_{p,q}^\alpha} \leq \frac{\epsilon}{s^\alpha} s^{\alpha-d/p} s^{d/p} = \epsilon.$$

573 For a fixed $f^* \in \mathcal{F}$, let $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ be a set of noisy observations with $y_i = f^*(\mathbf{x}_i) + \epsilon_i$, $\epsilon_i \sim$
 574 $\text{SubGaussian}(0, \sigma^2 I)$. Further assume that \mathbf{x}_i are evenly distributed in $(0, 1)^d$ such that in all
 575 regions as defined in (16), the number of samples is $n_j := O(n/s^d)$. Using Le Cam's inequality, we
 576 get that in any region, any estimator θ satisfy

$$\sup_{f^* \in \mathcal{F}} \mathbb{E}_{\mathcal{D}} [\|\theta(\mathcal{D}) - f^*\|_j] \geq \frac{C_m \epsilon}{16s^\alpha}$$

577 as long as $(\frac{\epsilon}{s^\alpha})^2 \lesssim \frac{s^d}{n}$, where $\|\cdot\|_j := \frac{1}{n_j} \sum_{s(\mathbf{x}-\mathbf{j}) \in [0,1]^d} |f(\mathbf{x})|$ denotes the norm defined in the
 578 block indexed by \mathbf{j} , C_m is a constant that depends only on m . Choosing $s = O(n^{\frac{1}{2\alpha+d}})$, we get

$$\sup_{f^* \in \mathcal{F}} \mathbb{E}_{\mathcal{D}} [\|\theta(\mathcal{D}) - f^*\|_j] \geq n^{-\frac{\alpha}{2\alpha+d}}.$$

579 Observing $\frac{1}{n} \sum_{i=1}^n L(\hat{f}(\mathbf{x}_i)) \geq 0.5 \sum_{i=1}^n |f(\mathbf{x}_i) - f^*(\mathbf{x}_i)| \approx \frac{1}{s^d} \sum_{\mathbf{j} \in [s]^d} \|\hat{f} - f^*\|_j$ finishes the
 580 proof.

581 F Supporting theorem

Lemma 15. [Lemma 14 in Zhang and Wang [27]] For any $a \in \mathbb{R}^{\bar{M}}$, $0 < p' < p$, it holds that:

$$\|a\|_{p'}^{p'} \leq \bar{M}^{1-p'/p} \|a\|_p^{p'}.$$

Proposition 16 (Proposition 7 in Zhang and Wang [27]). Let $\alpha - d/p > 1, r > 0$. For
 any function in Besov space $f^* \in B_{p,q}^\alpha$ and any positive integer \bar{M} , there is an \bar{M} -sparse
 approximation using B-spline basis of order m satisfying $0 < \alpha < \min(m, m-1+1/p)$:
 $\tilde{f}_{\bar{M}} = \sum_{i=1}^{\bar{M}} a_{k_i, \mathbf{s}_i} M_{m, k_i, \mathbf{s}_i}$ for any positive integer \bar{M} such that the approximation error is bounded
 as $\|\tilde{f}_{\bar{M}} - f^*\|_r \lesssim \bar{M}^{-\alpha/d} \|f^*\|_{B_{p,q}^\alpha}$, and the coefficients satisfy

$$\|\{2^{k_i} a_{k_i, \mathbf{s}_i}\}_{k_i, \mathbf{s}_i}\|_p \lesssim \|f^*\|_{B_{p,q}^\alpha}.$$

582 **Lemma 17** (Lemma 11 in [27]). Let $M_{m,k,s}$ be the B-spline of order m with scale 2^{-k} in each
 583 dimension and position $\mathbf{s} \in \mathbb{R}^d$: $M_{m,k,s}(\mathbf{x}) := M_m(2^k(\mathbf{x} - \mathbf{s}))$, M_m is defined in (4). There exists
 584 a neural network with d -dimensional input and one output, with width $w_{d,m} = O(dm)$ and depth
 585 $L \lesssim \log(C_{13}/\epsilon)$ for some constant C_{13} that depends only on m and d , approximates the B spline
 586 basis function $M_{m,k,s}(\mathbf{x}) := M_m(2^k(\mathbf{x} - \mathbf{s}))$. This neural network, denoted as $\tilde{M}_{m,k,s}(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^d$,
 587 satisfy

- 588 • $|\tilde{M}_{m,k,s}(\mathbf{x}) - M_{m,k,s}(\mathbf{x})| \leq \epsilon$, if $0 \leq 2^k(x_i - s_i) \leq m+1, \forall i \in [d]$,
- 589 • $\tilde{M}_{m,k,s}(\mathbf{x}) = 0$, otherwise.
- 590 • The total square norm of the weights is bounded by $2^{2k/L} C_{14} dmL$ for some universal
 591 constant C_{14} .

592 **Proposition 18.** For any feedforward neural network f with width w and depth L with bias, there
593 exists a feedforward neural network f' with width $w' = w + 1$ and depth $L' = L$, such that for any
594 \mathbf{x} , $f(\mathbf{x}) = f'([\mathbf{x}^T, 1]^T)$

Proof. Proof by construction: let the weights in the ℓ -th layer in f be \mathbf{W}_ℓ , and the bias be \mathbf{b}_ℓ , and choose the weight in the corresponding layer in f' be

$$\mathbf{W}'_\ell = \begin{bmatrix} \tilde{\mathbf{W}}_\ell & \tilde{\mathbf{b}}_\ell \\ 0 & 1 \end{bmatrix}, \quad \forall \ell < L; \quad \mathbf{W}'_L = [\tilde{\mathbf{W}}_L \quad \tilde{\mathbf{b}}_L].$$

595 The constructed neural network gives the same output as the original one. \square

Corollary 19 (Corollary 13.7 and Corollary 14.3 in Wainwright [23]). Let

$$\mathcal{G}_n(\delta, \mathcal{F}) = \mathbb{E}_{w_i} \left[\sup_{g \in \mathcal{F}, \|g\|_n \leq \delta} \left| \frac{1}{n} \sum_{i=1}^n w_i g(\mathbf{x}_i) \right| \right], \quad \mathcal{R}_n(\delta, \mathcal{F}) = \mathbb{E}_{\epsilon_i} \left[\sup_{g \in \mathcal{F}, \|g\|_n \leq \delta} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i g(\mathbf{x}_i) \right| \right],$$

596 denotes the local Gaussian complexity and local Rademacher complexity respectively, where $w_i \sim$
597 $\mathcal{N}(0, 1)$ are the i.i.d. Gaussian random variables, and $\epsilon_i \sim \text{uniform}\{-1, 1\}$ are the Rademacher
598 random variables. Suppose that the function class \mathcal{F} is star-shaped, for any $\sigma > 0$, any $\delta \in (0, \sigma]$
599 such that

$$\frac{16}{\sqrt{n}} \int_{\delta_n^2/4\sigma}^{\delta_n} \sqrt{\log \mathcal{N}(\mathcal{F}, \mu, \|\cdot\|_\infty)} d\mu \leq \frac{\delta_n^2}{4\sigma}$$

600 satisfies

$$\mathcal{G}_n(\delta, \mathcal{F}) \leq \frac{\delta^2}{2\sigma}. \quad (17)$$

601 Furthermore, if \mathcal{F} is uniformly bounded by b , i.e. $\forall f \in \mathcal{F}, \mathbf{x} |f(\mathbf{x})| \leq b$ any $\delta > 0$ such that

$$\frac{64}{\sqrt{n}} \int_{\delta_n^2/2b4\sigma}^{\delta_n} \sqrt{\log \mathcal{N}(\mathcal{F}, \mu, \|\cdot\|_\infty)} d\mu \leq \frac{\delta_n^2}{b}.$$

602 satisfies

$$\mathcal{R}_n(\delta, \mathcal{F}) \leq \frac{\delta^2}{b}. \quad (18)$$

603 **Proposition 20.** An L -layer ReLU neural network with no bias and bounded norm

$$\sum_{\ell=1}^L \|\mathbf{W}_\ell\|_{\text{F}}^2 \leq B$$

604 is Lipschitz continuous with Lipschitz constant $(B/L)^{L/2}$

605 *Proof.* Notice that ReLU function is 1-homogeneous, similar to Proposition 4 in [27], for any neural
606 network there exists an equivalent model satisfying $\|\mathbf{W}_\ell\|_{\text{F}} = \|\mathbf{W}_{\ell'}\|_{\text{F}}$ for any ℓ, ℓ' , and its total
607 norm of parameters is no larger than the original model. Because of that, it suffices to consider the
608 neural network satisfying $\|\mathbf{W}_\ell\|_{\text{F}} \leq \sqrt{B/L}$ for all ℓ . The Lipschitz constant of such linear layer is
609 $\|\mathbf{W}_\ell\|_2 \leq \|\mathbf{W}_\ell\|_{\text{F}} \leq \sqrt{B/L}$, and the Lipschitz constant of ReLU layer is 1. Taking the product
610 over all layers finishes the proof. \square

611 **Proposition 21.** An L -layer ReLU convolution neural network with convolution kernel size K , no
612 bias and bounded norm

$$\sum_{\ell=1}^L \|\mathbf{W}_\ell\|_{\text{F}}^2 \leq B.$$

613 is Lipschitz continuous with Lipschitz constant $(KB/L)^{L/2}$

614 This proposition can be proved by taking Proposition 24 into the proof of Proposition 20.

615 **Proposition 22.** Let $f = f_{\text{post}} \circ (1 + f_{\text{NN}} + f_{\text{other}}) \circ f_{\text{pre}}$ be a ResNeXt, where $1 + f_{\text{NN}} + f_{\text{other}}$
616 denotes a residual block, f_{pre} and f_{post} denotes the part of the neural network before and after this
617 residual block, respectively. f_{NN} denotes one of the building block in this residual block and f_{other}
618 denotes the other residual blocks. Assume $f_{\text{pre}}, f_{\text{NN}}, f_{\text{post}}$ are Lipschitz continuous with Lipschitz
619 constant $L_{\text{pre}}, L_{\text{NN}}, L_{\text{post}}$ respectively. Let the input be \mathbf{x} , if the residual block is removed, the
620 perturbation to the output is no more than $L_{\text{pre}} L_{\text{NN}} L_{\text{post}} \|\mathbf{x}\|$

Proof.

$$\begin{aligned}
& |f_{\text{post}} \circ (1 + f_{\text{NN}} + f_{\text{other}}) \circ f_{\text{pre}}(\mathbf{x}) - f_{\text{post}} \circ (1 + f_{\text{other}}) \circ f_{\text{pre}}(\mathbf{x})| \\
& \leq L_{\text{post}} |(1 + f_{\text{NN}} + f_{\text{other}}) \circ f_{\text{pre}}(\mathbf{x}) - (1 + f_{\text{other}}) \circ f_{\text{pre}}(\mathbf{x})| \\
& = L_{\text{post}} |f_{\text{NN}} \circ f_{\text{pre}}(\mathbf{x})| \\
& \leq L_{\text{pre}} L_{\text{NN}} L_{\text{post}} \|\mathbf{x}\|.
\end{aligned}$$

621

□

622 **Proposition 23.** *The neural network defined in Lemma 6 with arbitrary number of blocks has*
623 *Lipschitz constant $\exp((KB_{\text{res}}/L)^{L/2})$, where $K = 1$ when the feedforward neural network is the*
624 *building blocks and K is the size of the convolution kernel when the convolution neural network is*
625 *the building blocks.*

626 *Proof.* Note that the m -th block in the neural network defined in Lemma 6 can be represented
627 as $y = f_m(\mathbf{x}; \omega_m) + \mathbf{x}$, where f_m is an L -layer feedforward neural network with no bias. By
628 Proposition 20 and Proposition 21, such block is Lipschitz continuous with Lipschitz constant
629 $1 + (KB_m/L)^{L/2}$, where the weight parameters of the m -th block satisfy that $\sum_{\ell=1}^L \|W_\ell^{(m)}\|_{\mathbb{F}}^2 \leq B_m$
630 and $\sum_{m=1}^M B_m \leq B_{\text{res}}$.

631 Since the neural network defined in Lemma 6 is a composition of M blocks, it is Lipschitz with
632 Lipschitz constant L_{res} . We have

$$L_{\text{res}} \leq \prod_{m=1}^M \left(1 + \left(\frac{KB_m}{L} \right)^{L/2} \right) \leq \exp \left(\sum_{m=1}^M \left(\frac{KB_m}{L} \right)^{L/2} \right),$$

633 where we use the inequality $1 + z \leq \exp(x)$ for any $x \in \mathbb{R}$. Furthermore, notice that
634 $\sum_{m=1}^M (KB_m/L)^{L/2}$ is convex with respect to (B_1, B_2, \dots, B_M) when $L > 2$. Since $\sum_{m=1}^M B_m \leq$
635 B_{res} and $B_m \geq 0$, then we have $\sum_{m=1}^M (KB_m/L)^{L/2} \leq (KB_{\text{res}}/L)^{L/2}$ by convexity. Therefore,
636 we obtain that $L_{\text{res}} \leq \exp((KB_{\text{res}}/L)^{L/2})$. □

637 **Proposition 24.** *For any $\mathbf{x} \in \mathbb{R}^d$, $\mathbf{w} \in \mathbb{R}^K$, $K \leq d$, $\|\text{Conv}(\mathbf{x}, \mathbf{w})\|_2 \leq \sqrt{K} \|\mathbf{x}\|_2 \|\mathbf{w}\|_2$.*

638 *Proof.* For simplicity, denote $x_i = 0$ for $i \leq 0$ or $i > d$.

$$\begin{aligned}
\|\text{Conv}(\mathbf{x}, \mathbf{w})\|_2^2 &= \sum_{i=1}^d \langle \mathbf{x}[i - \frac{K-1}{2} : i + \frac{K-1}{2}], \mathbf{w} \rangle^2 \\
&\leq \sum_{i=1}^d \|\mathbf{x}[i - \frac{K-1}{2} : i + \frac{K-1}{2}]\|_2^2 \|\mathbf{w}\|_2^2 \\
&\leq K \|\mathbf{x}\|_2^2 \|\mathbf{w}\|_2^2,
\end{aligned}$$

639 where the second line comes from Cauchy-Schwarz inequality, the third line comes by expanding
640 $\|\mathbf{x}[i - \frac{K-1}{2} : i + \frac{K-1}{2}]\|_2^2$ by definition and observing that each element in \mathbf{x} appears at most K
641 times. □