

# On a Sparse Shortcut Topology of Artificial Neural Networks

Feng-Lei Fan<sup>1</sup>, *Member, IEEE*, Dayang Wang<sup>2</sup>, Hengtao Guo<sup>1</sup>, Qikui Zhu<sup>1</sup>, *Member, IEEE*,  
Pingkun Yan<sup>1\*</sup>, *Senior Member, IEEE*, Ge Wang<sup>1\*</sup>, *Fellow, IEEE*, and Hengyong Yu<sup>2\*</sup>, *Senior Member, IEEE*

**Abstract**—In established network architectures, shortcut connections are often used to take the outputs of earlier layers as additional inputs to later layers. Despite the extraordinary effectiveness of shortcuts, there remain open questions on the mechanism and characteristics. For example, why are shortcuts powerful? Why do shortcuts generalize well? In this paper, we investigate the expressivity and generalizability of a novel sparse shortcut topology. First, we demonstrate that this topology can empower a one-neuron-wide deep network to approximate any univariate continuous function. Then, we present a novel width-bounded universal approximator in contrast to depth-bounded universal approximators and extend the approximation result to a family of equally competent networks. Furthermore, with generalization bound theory, we show that the proposed shortcut topology enjoys excellent generalizability. Finally, we corroborate our theoretical analyses by comparing the proposed topology with popular architectures, including ResNet and DenseNet, on well-known benchmarks and perform a saliency map analysis to interpret the proposed topology. Our work helps enhance the understanding of the role of shortcuts and suggests further opportunities to innovate neural architectures.

**Impact Statement**—Shortcuts are the key elements of many well-performed neural network architectures and have achieved huge success in many applications. However, over the past years, why shortcuts are powerful was not so much investigated from a theoretical point of view. To fill this gap, we present detailed analyses on the power of a sparse shortcut topology in views of expressivity and generalizability. Furthermore, our theoretical studies are corroborated by comprehensive prediction and classification experiments. Our work is useful in understanding the role of shortcuts and can inspire more research in neural architecture design.

**Index Terms**—Theoretical deep learning, network architecture, shortcut network, expressivity, generalizability

## I. INTRODUCTION

Recently, deep learning [1] has been rapidly evolving and achieved great success in many applications [2]–[6]. Since AlexNet [7], more and more models were developed; for example, Inception [8], Network in Network [9], VGG [10], ResNet [11], DenseNet [12], and so on. These models play an important role as backbone architectures, pushing the performance boundaries of deep learning on the downstream tasks.

\*Drs. Pingkun Yan, Ge Wang and Hengyong Yu serve as co-corresponding authors. This work was supported by IBM AI Horizon Scholarship, R01EB026646, R01CA233888, R01CA237267, R01HL151561, R21CA264772, and R01EB031102.

<sup>1</sup>Feng-Lei Fan (fanf2@rpi.edu), Hengtao Guo, Qikui Zhu, Pingkun Yan and Ge Wang (wangg6@rpi.edu) are with Department of Biomedical Engineering, Rensselaer Polytechnic Institute, Troy, NY 12180, USA

<sup>2</sup>Dayang Wang and Hengyong Yu (hengyong\_yu@uml.edu) are with Department of Electrical and Computer Engineering, University of Massachusetts, Lowell, MA 01854, USA

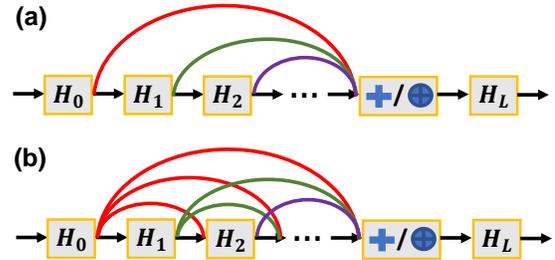


Fig. 1. Comparison of sparse and dense shortcut topologies. (a) A novel sparse shortcut topology; (b) the densely connected topology, where  $H_i$  denotes a collection of common operations such as convolution, ReLU, and so on. There are two aggregation methods: summation and concatenation marked as  $+$  and  $\oplus$ , respectively. In this paper, the summation is used for expressivity and concatenation for generalization purposes.

In these studies, great efforts were made to explore the use of skip connections [13]–[17]. For instance, a shortcut topology was searched in the framework of a lightweight network for a super-resolution task [13]. Hypercolumn Network [14] stacked the units at all layers as a concatenated feature descriptor to obtain semantic information and precise localization. Highway Network [16] achieved great success in training a very deep network. Fractal Network [17] utilized a different skip connection design, by which interacting sub-paths were used without any pass-through or residual connections.

In the 1990s, the universal approximation theorem was proved to justify the representation power of a network. Given a sufficient number of neurons, a one-hidden-layer network can express any continuous function [18], [19]. Recently, inspired by the success of deep learning, intensive efforts were put to explain the advantages of depth over width of a network. The basic idea behind these results is to construct a particular class of functions that a deep network can efficiently represent, but shallow networks cannot [20]–[24]. However, despite incorporating shortcuts greatly empowers a neural network in solving real-world problems, theoretical studies are few to explain the representation and generalization abilities of shortcuts. In this study, we present our theoretical findings on a novel sparse shortcut topology, wherein shortcuts are used to bridge all prior layers and the final layer in a block or the whole network (see Figure 1(a)), thereby partially addressing why shortcuts are effective types of machinery in a network.

First, we show that a one-neuron-wide network with the proposed topology can approximate any univariate function, while a one-neuron-wide feedforward network cannot. This suggests that adding shortcuts can lead to a more powerful network structure. Along this direction, we report an alterna-

tive novel width-bounded universal approximator by using the Kolmogorov-Arnold representation theorem [25], in contrast to the depth-bounded universal approximator [26]–[28]. The width-bounded universal approximator refers to the universal approximators whose width is limited, but depth is arbitrarily large, while the depth-bounded universal approximator has a limited depth, but its width can be arbitrarily large. Given the input of  $n$  dimensions, the required width is no more than  $2n^2+n$  per layer in our scheme. Then, we extend the result to a family of networks such that given approximation ability, these networks are equally competent. Furthermore, we analyze the effect of concatenation shortcuts on the generalization bound of deep networks. We show that the investigated topology enjoys a tighter generalization bound compared with the densely connected one, which suggests that the investigated topology can generalize well. To verify the positive results from the theoretical analyses, we prototype a network with the proposed topology and evaluate its performance on some well-known benchmarks. Finally, the experimental results demonstrate that the constructed network can achieve competitive learning performance compared to networks with residual topologies, the densely connected network, and other state-of-the-art models.

In summary, our contributions are three-fold. 1) We demonstrate the expressivity of the shortcut connections by presenting a univariate continuous function approximation theorem and a width-limited universal approximator, which partially addresses why networks with shortcuts are powerful. 2) To the best of our knowledge, our work is the first to analyze the generalizability of concatenation shortcuts based on the generalization bound theory. In addition, we also show that the generalization bounds of the proposed topology are tighter than those of the densely connected topology. 3) We conduct experiments to validate our theoretical analyses, and the investigated topology performs competitively in regression and classification experiments on several well-known benchmarks.

To clarify, all our studies are based on the architecture shown in Figure 1(a), which is a construction of skip connections. The central hypothesis of this paper is that the proposed topology in Figure 1(a) enjoys good expressivity (Section III) and generalizability (Section IV). Because the core of the proposed topology is the employment of shortcuts, our work also explains why shortcuts are essential in a network structure. This hypothesis is validated by comprehensive experimental comparisons (Section V).

## II. RELATED WORK

There are studies to explain the success of summation shortcuts. It was reported in [28] that with residual connections, one neuron is sufficient for the ResNet to approximate any Lebesgue-integrable function. In [29], it was showcased that the residual networks demonstrate an ensemble-like behavior. Liu *et al.* [30] studied the convergence behavior of a two-layer network and proved that the optimization of a two-layer ResNet can avoid spurious minima under mild restrictions. He *et al.* [31] studied a spectrally-normalized margin bound to discuss the influence of residual connections on the generalization ability of deep networks. They showed that the

margin-based multi-class generalization bound of ResNet is of the same magnitude as that of chain-like counterparts. Therefore, the generalizability of ResNet is not worse than that of a feedforward network. Here, we not only justify the representation ability of summation shortcuts but also conduct the generalization bound analysis for concatenation shortcuts, which systematically enrich our understanding of the expressivity and generalizability of shortcuts.

The work closely related to ours was done in [32], [33], which utilized the proposed network topology (Figure 1(a)) as a backbone for CT image denoising and super-resolution. However, their studies were not theoretical and did not answer why such a structure can work. In contrast, we approach the utility of this shortcut topology through detailed mathematical analyses and comprehensive experiments. In addition, the investigated topology here is a sparsified version of the densely connected shortcut topology. By setting the relevant weights as zero, the densely connected topology will reduce into the topology here. Our results somehow show that the densely connected topology is redundant.

As far as the universal approximation is concerned, in Lu *et al.* [26], giving at most  $n+4$  neurons per layer and allowing an infinite depth, a fully-connected deep network with ReLU activation functions can accurately approximate a Lebesgue-integrable  $n$ -dimension function in the  $L^1$ -norm sense. As an extension, Lin *et al.* [28] compressed  $n+4$  into 1 by using residual connections. They also argued that because the identity mapping should be counted as  $n$  units, the actual width of their network is  $n+1$ . Along this direction, we exploit the Kolmogorov-Arnold representation theorem [25] to derive a novel width-limited universal approximator with a width no more than  $2n^2+n$  per layer. Although the upper bound of width in our universal approximator is greater than those set by [26] and [28], our work is still valuable because of the methodology novelty and the scarcity of width-bounded universal approximators.

## III. EXPRESSIVITY

In this section, we first study the representation ability of the shortcut topology shown in Figure 1(a) that is based on summation (+) aggregation by presenting its superior approximation ability and then extend the results to more shortcut topologies, thereby shedding light on the question why shortcuts are powerful.

### A. Univariate continuous function approximation

Our main result is that adding shortcuts, as shown in Figure 1(a), can make a one-neuron wide network approximate any univariate continuous function in the sense of the  $L^\infty$  distance. It should be pointed out that our result is constructive, and it is still an open problem to prove that the trained network converges to our construction. Mathematically, we make the following proposition:

**Proposition 1:** With ReLU activation functions for all hidden neurons, for any continuous function  $g : [0, 1] \rightarrow \mathbb{R}$  and any given precision  $\delta > 0$ , there exists a neural network

$G$  of the proposed topology with one neuron in each layer such that

$$\sup_{x \in [0,1]} |g(x) - G(x)| < \delta \quad (1)$$

**The sketch of our constructive analysis:** Any univariate continuous function can be approximated by a continuous piecewise linear function within any given closeness [34]. Therefore, the key of proof becomes how to implement this piecewise approximation by a one-neuron-wide network of the proposed topology. In our scheme, we use the ReLU as activation functions for all neurons except the output neuron. By the convention of regression tasks, the activation function of the output layer is linear. Our construction is to make each neuron represent a piecewise function, and then we use shortcuts to aggregate these piecewise linear segments in the output neuron.

**Preliminaries:** Without loss of generality, a continuous function  $g(x)$  can be approximated by a continuous piecewise linear function  $f(x)$  at any accuracy in the  $L^\infty$  sense, provided that the interval  $[0, 1]$  is partitioned into very tiny sub-intervals. Therefore, to demonstrate the correctness of **Proposition 1**, we just need to use a one-neuron-wide network of the investigated topology to implement  $f(x)$ . Suppose that there are  $N$  pieces in  $f(x)$ , we can construct an explicit expression of  $f(x)$  as follows:

$$f(x) = \begin{cases} f_0(x) & x \in [x_0, x_1] \\ f_1(x) & x \in (x_1, x_2] \\ \vdots & \\ f_{N-1}(x) & x \in (x_{N-1}, x_N] \end{cases}, \quad (2)$$

where  $x_0 = 0$ ,  $x_N = 1$ , and

$$f_i(x) = \begin{cases} \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i} (x - x_i) + f(x_i) & x \in [x_i, x_{i+1}] \\ 0 & x \notin [x_i, x_{i+1}] \end{cases} \quad (3)$$

for  $i = 0, 1, 2, \dots, N-1$ , satisfying continuity. Hereafter, we use  $M_i = \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i}$  for simplicity. By default, neighboring segments should have different slopes; otherwise they will be combined as one segment.

**Analysis:** Now, let us show how to select parameters of a one-neuron-wide network to express  $f(x)$  in the form of Eq. (2). The outputs of neurons are respectively denoted as  $R_0, R_1, R_2, \dots, R_{N-1}$ . For the  $i^{\text{th}}$  neuron, its output  $R_i$  is expressed as

$$R_i = (W_i x + b_i)^+, \quad (4)$$

where  $(\cdot)^+$  denotes the ReLU operation,  $W_i$  and  $b_i$  are the weight and bias respectively. In the following, mathematical induction is used to show that our construction can express  $f(x)$  exactly.

**Initial Condition  $R_0$ :** We use  $R_0$  to implement the linear function in the first interval  $[x_0, x_1]$ . By setting  $W_0 = |M_0|$ ,  $b_0 = -|M_0|x_0$ , the specific function of the first neuron becomes  $R_0 = (|M_0|(x - x_0))^+$ , where the ReLU keeps the linearity when  $x > x_0$ .

**Recurrent Relation:** Suppose that we have obtained the desired  $i^{\text{th}}$  neuron  $R_i$ , we can proceed to design the  $(i+1)^{\text{th}}$

neuron with the goal of expressing the function  $|f_{i+1}(x) - f_{i+1}(x_{i+1})|$ , which is  $|f_{i+1}(x)|$  over the interval  $(x_{i+1}, x_{i+2}]$  without a constant lift. The tricky point is that the current neuron basically takes in the output of the previous neuron as the input, which is in the functional range instead of the input domain. Therefore, we need to perform an inverse affine transform:

$$R_{i+1} = \left( |M_{i+1} - M_i| \times \left( \frac{1}{|M_i - M_{i-1}|} R_i - x_{i+1} + x_i \right) \right)^+ \quad (5)$$

For notation completeness,  $M_{-1} = 0$ . The trick we use is to invert  $R_i$  back to the input domain and set the new slope as  $|M_{i+1} - M_i|$ , which cancels the effect of  $R_i$  imposed on  $x > x_{i+1}$ , equivalently limiting  $R_i$  to only work over  $(x_i, x_{i+1}]$  once  $R_i$  and  $R_{i+1}$  are added together. The parameters in the  $(i+1)^{\text{th}}$  module are chosen as follows:  $W_{i+1} = \frac{|M_{i+1} - M_i|}{|M_i - M_{i-1}|}$  and  $b_{i+1} = (-x_{i+1} + x_i)|M_{i+1} - M_i|$ .

Thanks to the recurrent relation, we can compute each  $R_i$  as  $(|M_i - M_{i-1}|(x - x_i))^+$ . We aggregate the outputs of those  $N$  pieces in the final neuron through shortcut connections to get the neural network  $G(x)$  as follows:

$$G(x) = \sum_{i=0}^{N-1} \text{sgn}(i) R_i + f(x_0), \quad (6)$$

wherein  $\text{sgn}(i) = 1$  when  $M_i - M_{i-1} > 0$  and  $\text{sgn}(i) = -1$  when  $M_i - M_{i-1} < 0$ . Because  $R_i(x) = (|M_i - M_{i-1}|(x - x_i))^+$ , for any  $x \in [x_k, x_{k+1}]$ ,

$$\begin{aligned} G(x) &= \sum_{i=0}^{N-1} \text{sgn}(i) R_i + f(x_0) \\ &= \sum_{i=0}^{N-1} \text{sgn}(i) (|M_i - M_{i-1}|(x - x_i))^+ + f(x_0) \\ &= \sum_{i=0}^{N-1} (M_i - M_{i-1})(x - x_i)^+ + f(x_0) \\ &= \sum_{i=0}^k (M_i - M_{i-1})(x - x_i) + f(x_0) \\ &= \sum_{i=0}^k (M_i - M_{i-1})x - \sum_{i=0}^k (M_i - M_{i-1})x_i + f(x_0) \\ &= M_k x - M_k x_k + \sum_{i=0}^{k-1} M_i (x_{i+1} - x_i) + f(x_0) \\ &= M_k (x - x_k) + \sum_{i=0}^{k-1} (f(x_{i+1}) - f(x_i)) + f(x_0) \\ &= M_k (x - x_k) + f(x_k) \\ &= f_k(x), \end{aligned} \quad (7)$$

which indicates that  $G(x)$  can exactly express  $f(x)$  in Eq. (2).

To illustrate the idea clearly, we exemplify  $\sum_{i=0}^2 \text{sgn}(i) R_i + f(x_0)$  as  $R_0 + R_1 - R_2 + f(x_0)$ , as shown in Figure 2.

Based on the above derivation, for any  $f(x)$  consisting of  $N$  piecewise linear segments, there will be a function

$f(x_0) + \sum_{i=0}^{N-1} R_i$  constructed by a one-neuron-wide  $N$ -layer network in the proposed topology that can exactly represent  $f(x)$ . Because  $f(x)$  can approximate any continuous univariate function, **Proposition 1** is verified.

Now, let us analyze the limit of  $N$ . Suppose  $g \in C^1 : [0, 1] \rightarrow \mathbb{R}$ , because  $|g(x) - g(y)| \leq \int_{|x-y| \leq \eta} |g'(s)| ds \leq \eta \|g'\|_\infty$ , where  $\|g'\|_\infty$  is the maximum absolute value of the derivative of  $g$ , a continuous piecewise linear function  $f$  can represent  $g$ :  $\sup_x |g - f| < \delta$ , as long as we partition  $[0, 1]$  into intervals whose lengths are smaller than  $\delta / \|g'\|_\infty$ . As a result, the required number of pieces is  $1 / (\delta / \|g'\|_\infty) = \|g'\|_\infty / \delta$ , and the needed neuron number  $N$  for  $G$  is also  $\|g'\|_\infty / \delta$ .

**Remark 1:** An exciting question is whether the densely connected topology in the DenseNet is necessary or not. Zhu *et al.* [35] experimentally demonstrated that a sparse version of DenseNet has been excellent in image classification. In contrast, our **Proposition 1** theoretically confirms that given the sufficient depth, the densely connected topology has certain redundancy given representation ability, since a one-neuron-wide network with the proposed topology can already work for general approximation.

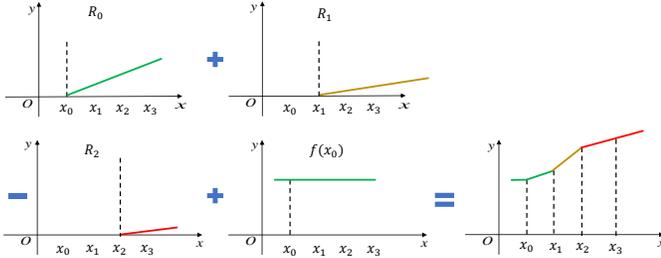


Fig. 2. An example of  $\sum_{i=0}^2 \text{sgn}(i)R_i + f(x_0)$  as  $R_0 + R_1 - R_2 + f(x_0)$  to illustrate how a one-neuron wide network can represent  $f(x)$ .

### B. Width-bounded universal approximator

Inspired by the feasibility of using a one-neuron-wide network to approximate any continuous univariate function, here we present an alternative width-bounded universal approximator, in analogy to a depth-bounded universal approximator. Width-bounded networks mean that the width of a network is limited, but the network can be arbitrarily deep. Our scheme is based on the topology in Figure 1(a) and the Kolmogorov-Arnold representation theorem. Specifically, we employ the Kolmogorov-Arnold representation theorem to bridge the gap between approximating univariate and multivariate functions.

**Proposition 2:** With ReLU activation functions, for any continuous function  $f : [0, 1]^n \rightarrow \mathbb{R}$  and any given precision  $\sigma > 0$ , there exists a neural network  $W$  with width no more than  $2n^2 + n$  per layer such that

$$\sup_{x_1, x_2, \dots, x_n \in [0, 1]} |f(x_1, x_2, \dots, x_n) - W(x_1, x_2, \dots, x_n)| < \sigma. \quad (8)$$

**Kolmogorov-Arnold representation theorem** [25]: For any continuous function  $f(x_1, \dots, x_n)$  with  $n \geq 2$ , there exist a group of continuous functions:  $\phi_{q,p}, q = 0, 1, \dots, 2n; p = 1, 2, \dots, n$  and  $\Phi_q$  such that

$$f(x_1, x_2, \dots, x_n) = \sum_{q=0}^{2n} \Phi_q \left( \sum_{p=1}^n \phi_{q,p}(x_p) \right). \quad (9)$$

**Scheme of analysis:** The representation theorem implies that any continuous function  $f(x_1, \dots, x_n)$  can be written as a composition of finitely many univariate functions. As shown in Figure 3, our scheme of approximating a multivariate continuous function  $f(x_1, \dots, x_n)$  is to first employ  $2n^2 + n$  single-neuron-wide sub-networks in the proposed topology to represent  $\phi_{q,p}(x_p)$  in a parallel manner. Next, suggested by the right side of Eq. (9), we summate the group of functions  $\{\phi_{q,1}(x_1), \phi_{q,2}(x_2), \dots, \phi_{q,n}(x_n)\}$  and feed  $\sum_{p=1}^n \phi_{q,p}(x_p)$  into a new one-neuron-wide network whose purpose is to approximate  $\Phi_q$ . Finally, we summate the yields of those  $2n + 1$  sub-networks as the ultimate output of the overall network.

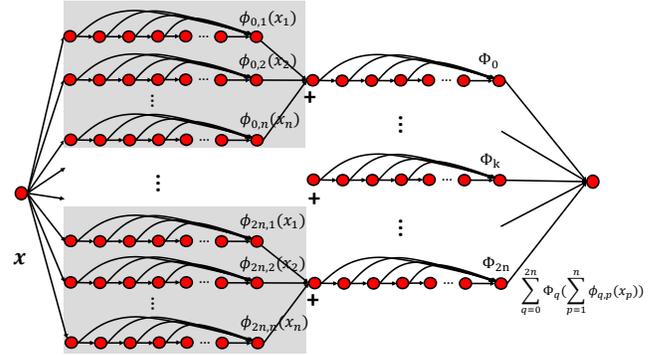


Fig. 3. The scheme of a width-bounded universal approximator.

**Analysis:** As we have shown in the **Proposition 1**, for every function  $\phi_{q,p}(x_p)$ , there exists a function  $D_{q,p}(x_p)$  represented by a one-neuron-wide network in the proposed topology such that

$$\sup_{x_p \in [0, 1]} |\phi_{q,p}(x_p) - D_{q,p}(x_p)| < \delta_{q,p}, \quad (10)$$

where  $\delta_{q,p} > 0$  is a given arbitrarily small quantity. After we integrate  $\{\phi_{q,1}(x_1), \phi_{q,2}(x_2), \dots, \phi_{q,n}(x_n)\}$ , for any selection of  $x_1, x_2, \dots, x_n \in [0, 1]$ , applying triangle inequality, we obtain the error of adding  $D_{q,p}$  with respect to  $p$  from Eq. (10):

$$\begin{aligned} & \sup_{x_1, x_2, \dots, x_n \in [0, 1]} \left| \sum_{p=1}^n \phi_{q,p}(x_p) - \sum_{p=1}^n D_{q,p}(x_p) \right| \\ & \leq \sup_{x_1, x_2, \dots, x_n \in [0, 1]} \sum_{p=1}^n |\phi_{q,p}(x_p) - D_{q,p}(x_p)| \\ & < \sum_{p=1}^n \delta_{q,p}. \end{aligned} \quad (11)$$

Given that  $\Phi_q$  is continuous, we employ the  $\epsilon - \delta$  definition of continuity: if  $g(x)$  is continuous at  $x_0$ , for any positive number  $\epsilon$ , there exists  $\delta(\epsilon, g) > 0$  satisfying that  $|g(x) - g(x_0)| < \epsilon$  when  $|x - x_0| < \delta$ . Let  $\epsilon = \frac{\sigma}{4n+2}$ , correspondingly we appropriately choose  $\delta_{q,p}$  so that  $\sum_{p=1}^n \delta_{q,p} < \delta(\frac{\sigma}{4n+2}, \Phi_q)$ . Thus, for every  $\Phi_q$ , we have the following:

$$\begin{aligned} & \sup_{x_1, x_2, \dots, x_n \in [0, 1]} \left| \Phi_q \left( \sum_{p=1}^n \phi_{q,p}(x_p) \right) - \Phi_q \left( \sum_{p=1}^n D_{q,p}(x_p) \right) \right| \\ & < \frac{\sigma}{4n+2}. \end{aligned} \quad (12)$$

Every continuous function  $\Phi_q$  is supported on  $\mathbb{R}$  instead of  $[0, 1]$ . Without loss of generality, we can still find a one-neuron-wide network in the proposed topology to approximate  $\Phi_q$  arbitrarily well. Let  $D_q(x)$  be the function expressed by such a network that can approximate  $\Phi_q$  in the precision of  $\frac{\sigma}{4n+2}$ , we have

$$\sup_{x \in \mathbb{R}} |\Phi_q(x) - D_q(x)| < \frac{\sigma}{4n+2}. \quad (13)$$

The above equation means that  $D_q(x)$  can represent  $\Phi_q(x)$  with an error no greater than  $\frac{\sigma}{4n+2}$  over  $\mathbb{R}$ . Introducing an intermediate term  $\Phi_q(\sum_{p=1}^n D_{q,p}(x_p))$  and applying the triangle inequality to estimate the error of feeding the summation of  $D_{q,p}$  into  $D_q$ , we have

$$\begin{aligned} & \sup_{x_1, x_2, \dots, x_n \in [0, 1]} \left| \Phi_q\left(\sum_{p=1}^n \phi_{q,p}(x_p)\right) - D_q\left(\sum_{p=1}^n D_{q,p}(x_p)\right) \right| \\ &= \sup_{x_1, x_2, \dots, x_n \in [0, 1]} \left| \Phi_q\left(\sum_{p=1}^n \phi_{q,p}(x_p)\right) - \Phi_q\left(\sum_{p=1}^n D_{q,p}(x_p)\right) \right| \\ &+ \left| \Phi_q\left(\sum_{p=1}^n D_{q,p}(x_p)\right) - D_q\left(\sum_{p=1}^n D_{q,p}(x_p)\right) \right| \\ &\leq \sup_{x_1, x_2, \dots, x_n \in [0, 1]} \left| \Phi_q\left(\sum_{p=1}^n \phi_{q,p}(x_p)\right) - \Phi_q\left(\sum_{p=1}^n D_{q,p}(x_p)\right) \right| \\ &+ \sup_{x_1, x_2, \dots, x_n \in [0, 1]} \left| \Phi_q\left(\sum_{p=1}^n D_{q,p}(x_p)\right) - D_q\left(\sum_{p=1}^n D_{q,p}(x_p)\right) \right| \\ &< \frac{\sigma}{4n+2} + \frac{\sigma}{4n+2} = \frac{\sigma}{2n+1}, \end{aligned} \quad (14)$$

where we enforce Eqs. (12) and (13) to derive from the second and third lines to the fourth line. Then, applying the triangle inequality for the summation of  $D_q, q = 0, \dots, 2n$ , we immediately obtain the error of the total approximation scheme of Kolmogorov-Arnold representation theorem from Eq. (9):

$$\begin{aligned} & \sup_{x_1, x_2, \dots, x_n \in [0, 1]} \left| \sum_{q=0}^{2n} \Phi_q\left(\sum_{p=1}^n \phi_{q,p}(x_p)\right) - \sum_{q=0}^{2n} D_q\left(\sum_{p=1}^n D_{q,p}(x_p)\right) \right| \\ &\leq \sup_{x_1, x_2, \dots, x_n \in [0, 1]} \sum_{q=0}^{2n} \left| \Phi_q\left(\sum_{p=1}^n \phi_{q,p}(x_p)\right) - D_q\left(\sum_{p=1}^n D_{q,p}(x_p)\right) \right| \\ &< (2n+1) \times \frac{\sigma}{2n+1} = \sigma. \end{aligned} \quad (15)$$

Let  $W(x_1, x_2, \dots, x_n) = \sum_{q=0}^{2n} D_q(\sum_{p=1}^n D_{q,p}(x_p))$ , we immediately get the validity of **Proposition 2**.

**Remark 2:** Here, we present a novel width-limited universal approximator with a width of no more than  $2n^2 + n$  per layer. This width bound is greater than those of other width-bounded universal approximators, e.g.,  $n+4$  in [26] and  $n+1$  in [28]. In addition, this bound is also greater than the width of common models. For example, the wide residual networks (WRN) have a width of 192, smaller than our bound. Despite that the width bound here is not pragmatic, due to the scarcity of width-bounded universal approximators and the novelty of our construction, it is still a valuable addition to the existing work. Moreover, the Kolmogorov-Arnold

representation theorem was revisited in [36]. The smoothness property of interior functions  $\phi_{q,p}$  of the Kolmogorov-Arnold representation was enhanced by modifying the interior functions as a mapping from digits of a binary expansion to digits of a ternary expansion. Such a modification enables a ReLU network to realize the modified Kolmogorov-Arnold representation. However, the resultant network has  $2K+3$  layers with  $\{n, 4n, \dots, 4n, n, 1, 2^{K^n}+1, 1\}$  neurons at each layer, respectively, where  $K$  is a positive number whose value is up to the pre-specified approximation precision. Such a network is neither depth-bounded nor width-bounded.

### C. A family of networks

Motivated by our constructive proof for the proposed topology, we report that in the one-dimensional setting, the aforementioned analysis is translatable to a rather inclusive family of network topologies. This network family (denoted as  $\Omega^M$ ) subsumes an extremely wide network, an extremely deep network, and networks between them, where  $M$  is the number of hidden neurons, not including the input and output nodes. We argue that network topologies in  $\Omega^M$  are equivalent in the sense of the approximation ability.

The input node is also considered as the neuron for simplicity. Hence, we refer to neurons as three types: hidden neurons, the input neuron, and the output neuron. A network in  $\Omega^M$  shall satisfy the following three conditions:

- 1) Every hidden neuron has one inbound edge.
- 2) Every hidden neuron and the input neuron have one outbound edge that links to the output neuron.
- 3) The input neuron is wired with at least one hidden neuron.

The first condition can be trivially relaxed to that every hidden neuron has multiple inbound edges by setting weights of extra edges as zero. The examples that belong to  $\Omega^6$  are shown in Figure 4. For a topology in  $\Omega^M$ , the number of required edges should be  $2M+1$ . One thing worthwhile to highlight is that members in  $\Omega^M$  are mutually convertible through one or more cutting-rewiring operations. A cutting-rewiring process means cutting the current input edge of one neuron and rewiring the one with another neuron. Regarding the network belonging to a network family  $\Omega$ , we have the following proposition:

**Proposition 3:** With ReLU activation functions, for any continuous function  $g : [0, 1] \rightarrow \mathbb{R}$  and any given precision  $\delta > 0$ , there is a network family  $\Omega^N$  in which any network  $K$ , whose mapping is denoted as  $\Omega_K^N(x)$ , satisfies:

$$\sup_{x \in [0, 1]} |g(x) - \Omega_K^N(x)| < \delta. \quad (16)$$

**The sketch of analysis:** Similarly, the core of the problem is how to represent a continuous piecewise function  $f(x)$  of  $N$  pieces by a network from  $\Omega^N(x)$ . The main difference is that the hidden neuron in a network from  $\Omega^N(x)$  is allowed to get the information from any previous neurons other than just precisely from the last neighboring neuron.

**Analysis:** For convenience and without loss of generality, we still use  $f(x)$  in Eq. (2). To prove **Proposition 3**, we need to use  $\Omega_K^N(x)$  to express  $f(x)$ .

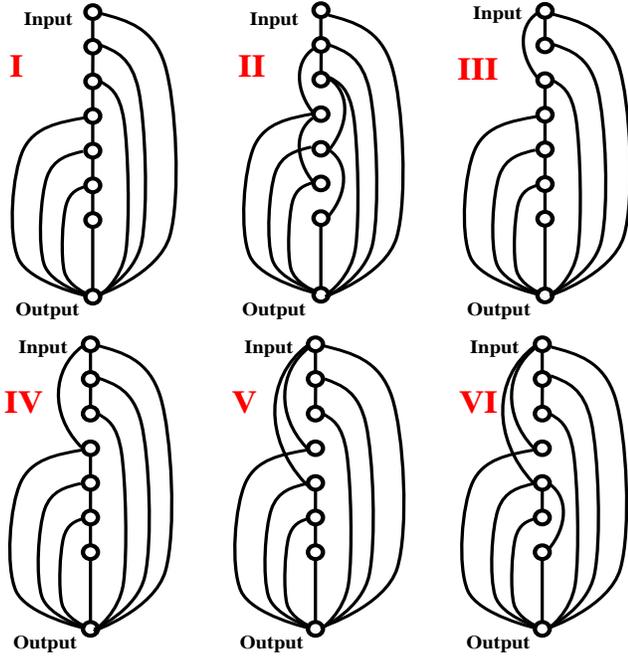


Fig. 4. Six exemplary structures in  $\Omega^6$  combined with ResNet setup are used to test if the networks in  $\Omega$  are truly equivalent or not.

Now we show how weights and bias in each neuron are appropriately selected in  $\Omega_K^N(x)$  to approximate  $f(x)$ . Without loss of generality, the neurons are denoted as  $Q_{input}, Q_0, \dots, Q_{N-1}, Q_{output}$ , where  $Q_{input}$  is the input node,  $Q_0$  is connected to the input neuron directly and  $Q_{i+1}$  is fed with either the input neuron or another neuron  $Q_t, t \leq i$ , and  $Q_{output}$  is the output neuron. Accordingly, the outputs of neurons  $Q_0, Q_1, \dots, Q_{N-1}$  are also denoted as  $Q_0, Q_1, \dots, Q_{N-1}$  for convenience, and our goal is to let  $Q_0, Q_1, \dots, Q_{N-1}$  to represent  $f_0, f_1, f_2, \dots, f_{N-1}$  at  $[x_0, x_1], (x_1, x_2), \dots, (x_{N-1}, x_N]$  without a constant shift.

For  $Q_0$ , similar to what we did before, we set that

$$Q_0 = (|M_0|(x - x_0))^+. \quad (17)$$

For  $Q_{i+1}$ , suppose that it connects with  $Q_j$ , we set

$$Q_{i+1} = \left( |M_{i+1} - M_i| \times \left( \frac{1}{|M_j - M_{j-1}|} Q_j - x_{i+1} + x_j \right) \right)^+. \quad (18)$$

Thus, the output of each neuron fulfills  $Q_i(x) = (|M_i - M_{i-1}|(x - x_i))^+$ . Similarly, we aggregate the output of all  $N$  hidden neurons in the output neuron as

$$\Omega_K^N(x) = \sum_{i=0}^{N-1} \text{sgn}(i) Q_i + f(x_0), \quad (19)$$

which is equal to  $f(x)$  according to Eq. (7). Therefore, we conclude **Proposition 3**.

**Remark 3:** Our representation ability analysis suggests that the members of  $\Omega$  are equivalently expressive. We want to emphasize that such a finding is important in both theoretical and practical senses. On the one hand, both a one-hidden-layer but super wide network and a one-neuron-wide but super deep network are demonstrated to have a strong expressive

ability. A natural curiosity is what about the networks in between. Do they also permit a good approximation ability? Here, we partially answer this question in the one-dimensional setting by showing that a wide network, a deep network, and networks in between from the network family  $\Omega$  are equally capable. On the other hand, network design is an important research direction. The insight can be drawn from our finding to network architecture design and search [13]. Since many networks are actually equivalent to each other, the search and design cost will be much reduced in principle.

#### IV. GENERALIZATION BOUND ANALYSIS

As mentioned earlier, for a shortcut topology, there are two types of aggregations: summation (+) and concatenation ( $\oplus$ ). The effect of summation connections on the generalizability of deep networks has been studied in [31]. To fill the gap that the effect of concatenation shortcuts is not explored, in this section, we dissect the generalizability of concatenation shortcuts by computing the generalization bounds. A generalization bound quantifying the generalization ability of a model is the upper bound of the generalization error. Recently, aimed at explaining good generalizability of over-parameterized deep networks, a plethora of norm-based generalization bounds [37]–[39] that rely on weight matrices norms rather than the number of weights have been developed. These bounds have a better explanation because they eliminate the direct dependence on the number of parameters.

Here, we derive the norm-based generalization bounds of DenseNet, with an emphasis on the spectrally normalized margin-based generalization bound [37]. To the best of our knowledge, our study is the first to analyze the effect of concatenation shortcuts on the generalization ability of deep networks. Then, we show that the generalization bound of the network using the proposed topology is tighter than that of the DenseNet, which suggests that the proposed topology can generalize well.

First, the data norm is set to the  $l_2$  norm and the operator norm set to the spectral norm  $\|\cdot\|_\sigma$  defined as  $\|A\|_\sigma = \sup_{\|Z\|_2 \leq 1} \|AZ\|_2$ . Furthermore,  $\|\cdot\|_{p,q}$  is the matrix  $(p, q)$ -norm defined as  $\|A\|_{p,q} = \left( \sum_{i=1}^n \|A_{:,i}\|_p^q \right)^{1/q}$ .

Next, we denote the model as  $F(\mathbf{x})$  and define the margin operator  $\mathcal{M} : \mathbb{R}^k \times \{1, 2, \dots, k\} \rightarrow \mathbb{R}$  for the  $k$ -class classification task as  $F(\mathbf{x})_z - \max_{j \neq z} F(\mathbf{x})_j$  for the  $z^{\text{th}}$  ground truth class, where  $z$  is the class index, and the ramp function is

$$l_\gamma(r) = \begin{cases} 0 & r < -\gamma \\ 1 + r/\gamma & -\gamma \leq r \leq 0 \\ 1 & r > 0, \end{cases} \quad (20)$$

where  $\gamma$  is the margin controlling the slope of  $l_\gamma(r)$ . Then, the empirical ramp loss over the dataset  $D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  is

$$\hat{\mathcal{R}}_\gamma(F) = \frac{1}{n} \sum_{i=1}^n (l_\gamma(-\mathcal{M}(F(\mathbf{x}_i), y_i))). \quad (21)$$

Minimizing the empirical ramp loss is equivalent to maximizing the margin of the predicted classes in the dataset. With all notations and definitions, we have the following theorem:

**Theorem 1:** Let us fix nonlinear activation functions  $\sigma_1, \dots, \sigma_L$ , where  $\sigma_i$  is  $\rho_i$ -Lipschitz and  $\sigma_i(0) = 0$ . Furthermore, let the margin  $\gamma > 0$ , spectral norm bounds  $(s_1, \dots, s_L)$ , data bound  $B$ , and matrices  $(2, 1)$ -norm bounds  $(b_1, \dots, b_L)$  be given. Then, with at least  $1 - \delta$  probability over  $N$  samples  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$  with  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $\sqrt{\sum_i \|\mathbf{x}_i\|_2^2} \leq B$  are drawn from identical and independent distribution, every DenseNet in  $F_{\mathcal{A}} : \mathbb{R}^d \rightarrow \mathbb{R}^k$  defined as

$$\begin{cases} G_0 = X^T \\ F_1 = A_1 X^T \\ G_i = \sigma_i(F_i) \\ F_{i+1} = A_{i+1} \oplus_{k=0}^i G_k \\ F_L = A_L \oplus_{k=0}^{L-1} G_k, \end{cases} \quad (22)$$

where  $X \in \mathbb{R}^{N \times d}$  collects all data samples  $\{\mathbf{x}_i\}_{i=1}^N$ ,  $\oplus$  is the matrix concatenation along the row direction,  $\oplus_{k=0}^i G_k = G_1 \oplus G_2 \cdots \oplus G_k = [G_0; G_1; \dots; G_k]$ ,  $A_i$  is of  $d_i \times n_i$  with  $n_i = \sum_{k=0}^{i-1} d_k$ , the matrices  $\mathcal{A} = (A_1, \dots, A_L)$  with  $A_i \in \mathbb{R}^{d_i \times n_i}$ ,  $n_i = \sum_{k=0}^{i-1} d_k$  obey that  $\|A_i\|_{\sigma} \leq s_i$  and  $\|A_i^T\|_{2,1} \leq b_i$ , and  $L$  is the number of layers, satisfies

$$\begin{aligned} & Pr\{\arg \max_i F_{\mathcal{A}}(\mathbf{x}_i) \neq y\} - \hat{\mathcal{R}}_{\gamma}(F_{\mathcal{A}}) \\ & \leq \frac{8}{n^{3/2}} + 3\sqrt{\frac{\ln(1/\delta)}{2n}} + \\ & \frac{36B \ln(n) \prod_{i=1}^L (1 + \rho_i s_i)}{\gamma n} \sqrt{\sum_{i=1}^L \frac{\rho_i^2 b_i^2}{(1 + \rho_i s_i)^2} \ln(2d_i n_i)}, \end{aligned} \quad (23)$$

where  $\ln(\cdot)$  is the natural logarithm. For conciseness, we put the proof of **Theorem 1** in Part A of supplementary materials.

**Remark 4:** Please note that our result is based on the proof in [37], and is the first to apply the results on the chain-like networks into the networks with concatenation shortcuts to evaluate the impact of concatenation shortcuts on the generalization bound of deep networks. As shown in Table I, we compare the bounds of the DenseNet and chain-like network. Incorporating dense concatenation shortcuts leads to a higher generalization bound than the chain-like network due to the increased matrix size  $n_i = \sum_{k=0}^{i-1} d_k > d_{max}$ . However, the bounds of the DenseNet and chain-like network are close when small weight matrices are used in each layer. This result partially explains why the DenseNet performs well in a small filter size because, in this situation, the concatenation shortcuts only moderately elevate the generalization bound.

TABLE I

THE GENERALIZATION BOUNDS FOR THE DENSENET AND CHAIN-LIKE NETWORK.  $d_{max}$  IS THE MAXIMUM WIDTH.

Models	Generalization Bound
DenseNet	$\mathcal{O}\left(\prod_{i=1}^L (1 + \rho_i s_i) \sqrt{\sum_{i=1}^L \frac{\rho_i^2 b_i^2}{(1 + \rho_i s_i)^2} \ln(2d_i n_i)}\right)$
Chain-like	$\mathcal{O}\left(\prod_{i=1}^L (\rho_i s_i) \sqrt{\sum_{i=1}^L \frac{b_i^2}{s_i^2} \ln(2d_{max}^2)}\right)$

**Proposition 4:** The margin-based multi-class generalization bound of the network in the proposed topology is tighter than that of the DenseNet.

**Insight:** The core of the derived bound in Eq. (23) is the third term of the right side, which is mainly dependent upon the spectral norm bound  $s_i$  and the matrix  $(2, 1)$ -norm bound  $b_i$  of weight matrices. Because by adding imaginary shortcuts (setting the extra weight matrices as zeros), the proposed topology becomes a particular case of the DenseNet, the spectral norm bounds and matrix  $(2, 1)$ -norm bounds of the proposed topology are no more than those of the DenseNet. Consequently, the spectrally normalized margin-based generalization bound of the network in the proposed topology is tighter than that of the DenseNet.

**Analysis:** Let us derive the margin-based multi-class generalization bound of the network in the proposed topology and compare it with that of the DenseNet. To discriminate them, in the following we use the superscript  $(S)$  for the parameters pertaining to the former and the superscript  $(D)$  to the latter. Then, Eq. (23) turns into

$$\begin{aligned} & Pr\{\arg \max_i F_{\mathcal{A}}^{(D)}(\mathbf{x}_i) \neq y\} - \hat{\mathcal{R}}_{\gamma}(F_{\mathcal{A}}^{(D)}) \\ & \leq \frac{8}{n^{3/2}} + 3\sqrt{\frac{\ln(1/\delta)}{2n}} + \\ & \frac{36B \ln(n) \prod_{i=1}^L (1 + \rho_i s_i^{(D)})}{\gamma n} \sqrt{\sum_{i=1}^L \frac{\rho_i^2 b_i^{(D)2}}{(1 + \rho_i s_i^{(D)})^2} \ln(2d_i n_i^{(D)})}. \end{aligned} \quad (24)$$

For a fair comparison, we set the output dimension of each layer in the network of the proposed topology to the same as that of the DenseNet. Also, we use  $d_i$  for both networks. Let  $A_i^{(S)}$  be of  $d_i \times n_i^{(S)}$ , where  $n_i^{(S)} = d_{i-1}$ ,  $i \leq L-1$ ,  $n_L^{(S)} = \sum_{i=1}^{L-1} d_i$  and  $X \in \mathbb{R}^{n \times d}$ . The computational structure of the network of the proposed topology is

$$\begin{cases} G_0^{(S)} = X^T \\ F_1^{(S)} = A_1^{(S)} X^T \\ G_i^{(S)} = \sigma_i(F_i^{(S)}) \\ F_{i+1}^{(S)} = A_{i+1}^{(S)} G_i^{(S)}, i \leq L-2 \\ F_L^{(S)} = A_L^{(S)} \oplus_{k=0}^{L-1} G_k^{(S)}. \end{cases} \quad (25)$$

Without changing the final output, we can rewrite the above structure by adding imaginary shortcuts and setting the extra weight matrices as zeros,

$$\begin{cases} G_0^{(S)} = X^T \\ F_1^{(S)} = A_1^{(S)} X^T \\ G_i^{(S)} = \sigma_i(F_i^{(S)}) \\ F_{i+1}^{(S)} = [A_{i+1}^{(S)}; \mathbf{0}^{d_{i+1} \times \sum_{k=0}^i d_k}] [G_i^{(S)}; \mathbf{0}^{d_i \times n}; \dots; \mathbf{0}^{d_0 \times n}] \\ F_L^{(S)} = A_L^{(S)} \oplus_{k=0}^{L-1} G_k^{(S)}, \end{cases} \quad (26)$$

where  $\mathbf{0}^{C_1 \times C_2}$  means the zero matrix of  $C_1 \times C_2$ . The network in the proposed topology is a special DenseNet with specific weight matrices as zeros. We can estimate the generalization bound for the above zero-padded network Eq. (26) by mim-

icking the generalization bound of DenseNet:

$$\begin{aligned}
& Pr\{\arg \max_i F_{\mathcal{A}}^{(S)}(\mathbf{x})_i \neq y\} - \hat{\mathcal{R}}_{\gamma}(F_{\mathcal{A}}^{(S)}) \\
& \leq \frac{8}{n^{3/2}} + 3\sqrt{\frac{\ln(1/\delta)}{2n}} + \\
& \frac{36B\ln(n) \prod_{i=1}^L (1 + \rho_i s_i^{(S)})}{\gamma n} \sqrt{\sum_{i=1}^L \frac{\rho_i^2 b_i^{(S)2}}{(1 + \rho_i s_i^{(S)})^2} \ln(2d_i n_i^{(D)})},
\end{aligned} \tag{27}$$

where  $n_i^{(D)}$  is used because the matrix size has been enlarged to the same to that of the DenseNet.

To verify **Proposition 4**, we need to compare the bounds of DenseNet and the proposed topology (Eq. (24) vs Eq. (27)). According to the definition of the spectral norm, we have

$$\begin{aligned}
& \|A_i^{(S)}\|_{\sigma} \\
& = \sup_{\|Z\|_2 \leq 1} \|A_i^{(S)} Z^{(S)}\|_2 \\
& = \sup_{\|Z\|_2 \leq 1} \|[A_i^{(S)}, \mathbf{0}^{d_i \times (n_i^{(D)} - n_i^{(S)})}]\| \|[Z; \mathbf{0}^{(n_i^{(D)} - n_i^{(S)}) \times n}]\|_2 \\
& \leq \sup_{\|Z\|_2 \leq 1} \|A_i^{(D)} Z^{(D)}\|_2 \\
& = \|A_i^{(D)}\|_{\sigma},
\end{aligned} \tag{28}$$

where zero padding is to make  $[A_i^{(S)}, \mathbf{0}]$  have the same size as that of  $A_i^{(D)}$ . Therefore, we derive that

$$s_i^{(S)} \leq s_i^{(D)}, i = 1, \dots, L. \tag{29}$$

In the same spirit, we can also derive that

$$b_i^{(S)} \leq b_i^{(D)}, i = 1, \dots, L. \tag{30}$$

Combining Eqs. (29) and (30), we have

$$\begin{aligned}
& \prod_{i=1}^L (1 + \rho_i s_i^{(S)}) \sqrt{\sum_{i=1}^L \frac{\rho_i^2 b_i^{(S)2}}{(1 + \rho_i^2 s_i^{(S)})^2} \ln(2d_i n_i^{(D)})} \leq \\
& \prod_{i=1}^L (1 + \rho_i s_i^{(D)}) \sqrt{\sum_{i=1}^L \frac{\rho_i^2 b_i^{(D)2}}{(1 + \rho_i s_i^{(D)})^2} \ln(2d_i n_i^{(D)})},
\end{aligned} \tag{31}$$

which has validated **Proposition 4**.

**Remark 5:** Our representation and generalization analyses suggest that DenseNet has certain redundancy in representation ability and a higher generalization bound. However, the redundant structure of DenseNet may facilitate the over-parameterization effect, which may cause optimization and generalization merits. For instance, regarding merits in optimization, stochastic gradient descent (SGD) can find the global minimum in shallow or deep networks in the setting of over-parameterization because there is a large set of global minimizers in an overly parameterized network [40]–[42]. Over-parameterization is also beneficial for generalization [43], [44]. Recently, the deep double descent phenomenon (When the model complexity increases, the generalization error goes down first and then up. However, as the model complexity keeps increasing and surpasses the so-called "interpolation

threshold", the generalization error starts going down) has been widely observed in many deep models [44]. In light of the double descent phenomenon, the complexity of DenseNet likely lies beyond the interpolation threshold.

## V. EXPERIMENTS

In this section, we conduct prediction and classification experiments on well-known benchmarks to evaluate the expressivity, generalizability, and interpretability of the proposed topology. The expressivity experiments use summation (+) shortcuts, while other experiments use concatenation ( $\oplus$ ) shortcuts. The competitive performance on prediction and classification tasks shows that the proposed topology is a desirable architecture, as suggested by encouraging theoretical analyses. In addition, we also demonstrate the superior interpretability of the investigated topology given the saliency map.

### A. Expressivity

We compare the expressivity of the proposed topology and residual topology in the infinite-width limit, where the gradient descent makes little change to the weights of a network. The training of a neural network with infinite width in each layer turns into a kernel ridge regression [45] process with the so-called neural tangent kernel (NTK [46]). When one fixes the type of activation functions, the neural tangent kernel of a neural network is only determined by the topology and the depth of the network [47]. Figure 5 shows the structures of the proposed network and a residual network that uses pre-activation features. In the proposed network, the output of each dense layer is connected to a layer before the final dense layer for summation. We denote the depth of both networks as  $K+2$ , where  $K$  is the number of residual blocks or the number of layers that constitute the proposed topology. Two networks are the same except for shortcut architectures.

Let samples of the training dataset be  $\{(\mathbf{x}_i, y_i)\}_{i=1}$ , where  $\mathbf{x}_i$  is the input and  $y_i$  is the output, and assume that  $f(\boldsymbol{\theta}, x)$  denotes the output of a neural network, where  $\boldsymbol{\theta}$  are all parameters, the  $(i, j)$ -entry of the NTK kernel  $\mathbf{H}^*$  [47] is defined by

$$ker(\mathbf{x}_i, \mathbf{x}_j) = \mathbb{E}_{\boldsymbol{\theta} \sim \Theta} \left\langle \frac{\partial f(\boldsymbol{\theta}, \mathbf{x}_i)}{\partial \boldsymbol{\theta}}, \frac{\partial f(\boldsymbol{\theta}, \mathbf{x}_j)}{\partial \boldsymbol{\theta}} \right\rangle. \tag{32}$$

The inference process is deterministic:

$$f(\mathbf{x}) = [ker(\mathbf{x}, \mathbf{x}_1), \dots, ker(\mathbf{x}, \mathbf{x}_n)] \cdot (\mathbf{H}^*)^{-1} \mathbf{y}. \tag{33}$$

Since the kernel in the inference process is only determined by the topology, depth, and the activation function, the comparison in the NTK domain can avoid the impact of other hyper-parameters such as the network width and learning rules (learning rate, batch size, optimizer, epoch number, and so on), which helps reveal the difference in the representation ability between two topologies.

We use the Boston house prices dataset [48] as a testbed that has 13 attributes including the average number of rooms, pupil-teacher ratio, and so on. The task is to predict the house price based on the attributes of a house. The dataset is randomly split into a training set (90%) and a test set (10%).

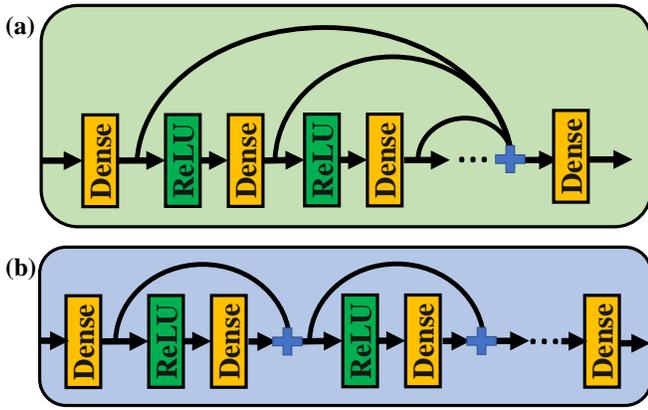


Fig. 5. "Dense" denotes a fully connected layer. (a) The network of the proposed topology; (b) the network of residual topology.

The mean squared error between predictions and ground truth is computed as the evaluation metric. We vary  $K$  from 4 to 10 to make a thorough comparison. The code is written online in Google Colab based on Python neural tangent package (<https://github.com/google/neural-tangents>). For all  $K$ , the inference time is no more than 10 seconds. Figure 6 highlights the consistent improvement of the proposed topology over the residual one. In addition, while the mean squared errors of both models keep going down as  $K$  increases, the downward momentum of the proposed topology is stronger.

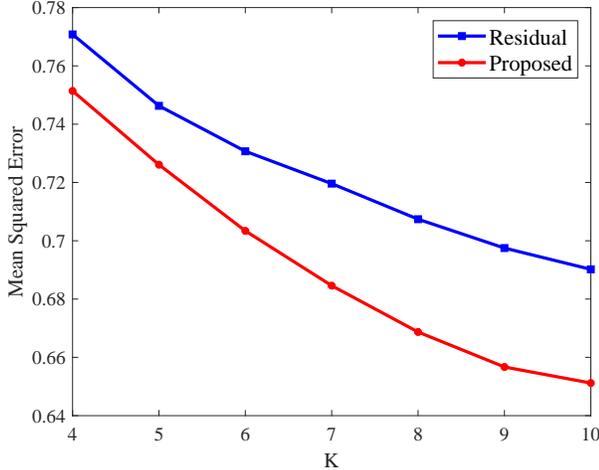


Fig. 6. Results of NTK kernel ridge regression of the residual topology and the proposed topology on the Boston house prices dataset.

### B. Generalizability

Here, we validate the generalizability of the proposed topology with concatenations to see if it can truly deliver competitive results as promised. Suppose that  $y_l$  is the output of the  $l^{\text{th}}$  module in the network of  $L$  modules, we characterize the workflow of the proposed topology in the following way:

$$\begin{aligned} y_{l+1} &= H_l(y_l), \\ y_L &= H_{L-1}(y_0 \oplus y_1 \oplus y_2 \oplus \dots \oplus y_{L-1}), \end{aligned} \quad (34)$$

where  $\oplus$  is a concatenation operator. The operator module  $H(\cdot)$  can perform multiple operations including batch normalization [49], convolution, dropout [50], and so on. While our theoretical analysis revolves around the multiplication operation, it can also scale to the convolution operation because a convolution between two vectors can be re-formulated as matrix multiplication.

The network of the proposed topology is implemented as a drop-in replacement for the DenseNet, which means that the only difference between our network and DenseNet is the shortcut topology. Our model comprises multiple blocks, and each block employs the proposed topology. Like DenseNet, the important hyperparameters for our model are the feature growth rate  $k$  and the number of layers in each block. The number of features of a layer in a block is referred to as the growth rate, which regulates the capacity of information passed to the final. We compare the proposed topology with other advanced deep learning benchmark models on the CIFAR-100, Tiny ImageNet, and ImageNet datasets.

**CIFAR-100:** We follow the initialization strategy in DenseNet. The DenseNet utilizes stage training: Across the stages, the number of filters is doubled, and the size of feature maps is reduced at the scale of 2. The proposed network includes four blocks. All model configurations for the proposed model follow the protocol in [12]. The total number of epochs is 250. The initial learning rate is 0.1 and divided by 10 in every quarter of the total epoch number. We use SGD for training with a weight decay of 0.0001 and a momentum of 0.9. We run each of the proposed models five times and compute the corresponding mean and variance of errors. In Table II, we summarize the experimental results on the CIFAR-100. The network of the proposed topology achieves a slightly higher error rate of 23.52% with much fewer parameters. The proposed model works better at larger growth rates, which is quite different from the DenseNet. Because of the memory constraint, a larger growth rate is prohibitive for the DenseNet. Overall, the proposed topology achieves competitive results over CIFAR-100.

TABLE II  
COMPARISONS OF TOP-1 ERRORS (%) ON CIFAR-100 AMONG THE PROPOSED MODEL AND OTHER MODELS.

Network	Params	Error(%)
NIN + Dropout [9]	-	35.68
FractalNet with Dropout [17]	38.6M	35.34
ResNet (Stocatic Depth) [51]	1.7M	37.80
DIANet [52]	-	23.02
SpinalNet [53]	-	35.01
LP-BNN [54]	-	23.02
DenseNet (k=12, depth=40)	1.0M	27.55
DenseNet (k=12, depth=100)	7.0M	23.79
DenseNet (k=24, depth=100)	27.2M	23.42
Proposed (k=12, depth=40)	0.4M	29.63 ± 0.017
Proposed (k=24, depth=40)	1.3M	26.21 ± 0.025
Proposed (k=40, depth=40)	3.6M	<b>23.52 ± 0.037</b>

The errors of compared models are reported by the official implementation.

**Tiny ImageNet dataset:** This dataset consists of 200 classes with 500 training, 50 validation, and 50 test images per class. The image size is  $64 \times 64$ , which are downsampled from the full images of the ImageNet dataset. In the experiments,

TABLE III  
COMPARISONS OF TOP-1 ERRORS (%) AMONG VARIOUS ADVANCING  
MODELS ON TINY IMAGENET.

Network	l.r.	Params	Error(%)
MobileNetV2 (2018) [55]	0.1	3.5M	43.76
EfficientNet-B0 (2019) [56]	0.1	5.3M	42.91
OctResNet50 (2019) [57]	0.1	25.5M	47.45
Lambda Network (2020) [58]	0.1	15.0M	58.71
SE-Net (2018) [59]	0.05	28.1M	53.98
Scale-Net (2019) [60]	0.01	31.4M	48.59
Ghost-Net (2020) [61]	0.1	5.2M	44.01
RandomWire-WS (2019) [62]	0.01	31.6M	42.11
Proposed A (k=96, depth=41)	0.1	5.0M	42.82± 0.31
<b>Proposed B (k=108, depth=41)</b>	0.1	10.6M	<b>42.04 ± 0.24</b>

All models are implemented by us.

we select the following models for comparison: MobileNetV2 [55], EfficientNet-B0 [56], OctResNet50 [57], Lambda Network [63], SE-Net [59], Scale-Net [60], Ghost-Net [61], and Randomly Wired Network [62]. All these models are well-known new benchmarks. We set the batch size to 64. We adopt the standard learning rate decay approach. In every 30 epochs, the learning rate is divided by 10. The initial learning rate is chosen from  $\{0.01, 0.05, 0.1\}$ . The momentum is 0.9. All models are trained in two TITAN Xp and one GeForce GTX 1080 GPUs. Among all models, it takes at most 742.43 seconds to finish one epoch. Based on our tuning, the appropriate hyperparameters for competitors are shown in Table III. We verify two models (k=96, depth=41, init-nf=32 and k=108, depth=41, init-nf=32), each of which consists of three blocks and "init-nf" means the number of features in the first layer of each block. We run the proposed two models five times and compute the mean and variance of errors. Table III shows top-1 validation errors of all models, where both proposed models achieve state-of-the-art performance. Particularly, the proposed model at a high growth rate obtains competitive accuracy over all the other models. One notable thing is that given a target performance, the network of the proposed topology uses three times fewer parameters than the randomly wired network.

**ImageNet dataset:** The ImageNet dataset [64] consists of 1.2 million images for training and 50,000 images for validation. No other augmentation techniques are employed in our experiments. We follow the basic data augmentation methods, as used in ShuffleNet, Randomly Wired Network, DenseNet, ECANet, and SENet. For model configurations, we follow those of DenseNet [12]. We set the batch size as 156, the initial learning rate as 0.1, the weight decay as 0.0001, and the momentum as 0.9. In validation, we adopt the standard 10-crop validation. To be fair, we compare our model with others in the small size regime ( $< 10M$  parameters) and regular size regime ( $\sim 20M$  parameters), respectively. It takes the smaller model around 75 minutes and the larger model around 90 minutes per epoch. We run the larger model three times and the smaller model five times to compute the average and variance of errors. Due to the computational burden of searches, NAS-based models appear in the small regime. Tables IV and V highlight the state-of-the-art results achieved by the proposed model. Regarding the small size

regime, despite a moderately higher model complexity, the proposed model achieves a performance superior or similar to those of other advanced models. Very favorably, our model is designed based on theoretical analyses. Compared to NAS, our model is free of computationally expensive searches. While for the regular model regime, our model is comparable to other advanced models.

TABLE IV  
THE TOP-1 ERROR (%) COMPARISONS IN SMALL MODEL REGIME ON  
IMAGENET VALIDATION SET.

Network	params	Error(%)
MobileNetV2 (2018) [55]	6.9M	25.3
ShuffleNet (2018) [65]	5.4M	26.3
NASNet-B (2018) [66]	5.3M	27.2
NASNet-C (2018) [66]	4.9M	27.5
Amoeba-A (2018) [67]	5.1M	25.5
Amoeba-B (2018) [67]	5.3M	26.0
PNAS (2018) [68]	5.1M	25.8
DARTS (2019) [69]	4.9M	26.9
FBNet-A (2019) [70]	4.3M	27.0
RandWire-WS (2019) [62]	5.6M	25.3 ± 0.25
RegNetX-600MF (2020) [71]	6.2M	25.9 ± 0.03
DeiT-Ti (2020) [72]	5.0M	25.4
<b>Proposed (k=96, depth=45)</b>	9.4M	<b>25.2 ± 0.07</b>

The errors of compared models are reported by the official implementation.

TABLE V  
THE TOP-1 ERROR (%) COMPARISONS IN REGULAR MODEL REGIME ON  
IMAGENET VALIDATION SET.

Network	params	Error(%)
SENet (2018) [59]	26.8M	23.3
ACNet (2019) [73]	19.8M	23.8
DenseNAS-R2 (2020) [74]	19.5M	24.2
ECA-Net (2020) [75]	24.4M	<b>22.5</b>
<b>Proposed (k=180, depth=41)</b>	27.9M	<b>22.9 ± 0.06</b>

The errors of compared models are reported by the official implementation.

### C. Interpretability

Interpretability is a fundamental problem for the development of deep learning [76], [77]. Here, we also show the superior interpretability of the proposed model in terms of the saliency map.

**Saliency map:** Currently, saliency methods deriving a saliency map by identifying relevance between features and the prediction of a model are the mainstream interpretability methods [78]. A myriad of saliency methods are based on gradients [76], the idea of which is that the strength of gradients can mirror the extent of how a feature can affect model output. As we know, shortcuts can facilitate training by alleviating gradient explosion and vanishing issues. The mechanism is that shortcuts provide additional paths for straightforward gradient propagation, improving the quality of saliency maps. In the proposed topology, shortcuts directly connect the final layer with all the prior layers, thereby conveying gradients among them. Meanwhile, the proposed topology can generalize, ignoring pixels from the input not located in the object that influences the image label, improving saliency maps as well. Integrating these two aspects, the saliency map of the proposed topology should be more accurate and sharper relative to the network without shortcuts.

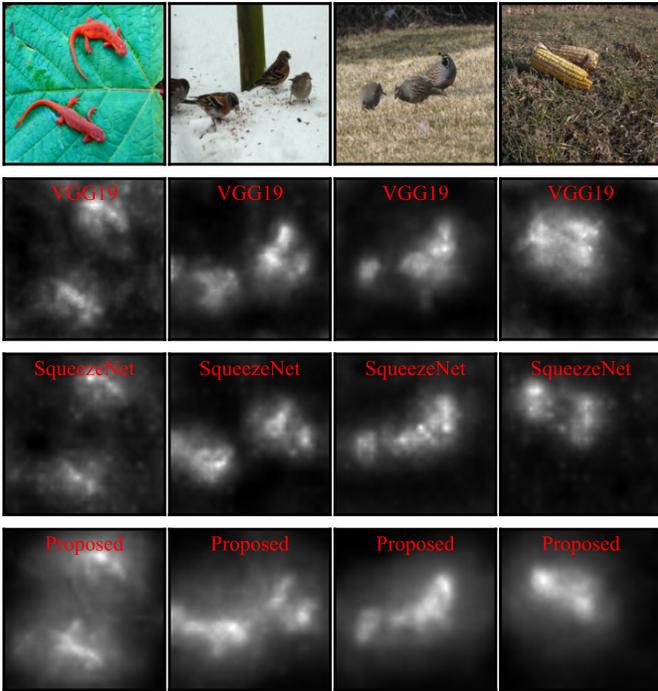


Fig. 7. Saliency maps of different models by the FullGrad method. Visually, regarding four images, saliency maps of the proposed model are sharper and their brightest points more conform to the objects.

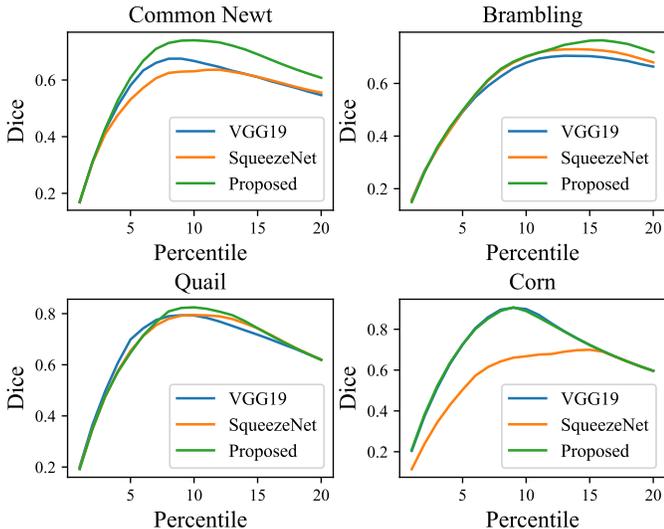


Fig. 8. Dice scores between the segmentation of an object and a saliency map as a function of the percentile. The segmentation of a saliency map is obtained by setting the  $q$ -percentile brightest pixels as one and the rest as zero.

We use the FullGrad method [79] to derive saliency maps because it can satisfy two characteristics (dependence and completeness) that the community has deemed important, while other classic methods such as SmoothGrad [80], IntegratedGrad [81], and so on cannot. Dependence describes that a feature is important if it can substantially affect the model output, while completeness is that the individual saliency scores must add up to the model output, which ensures that the total relevance corresponds to the extent of what is detected by a model. We compare our model with classic deep learning

models: VGG19 [10] and SqueezeNet [82]. Both models have no shortcuts. We have obtained three our models from three runs in the ImageNet experiments in the regular model regime. VGG19 and SqueezeNet are straightforwardly obtained from the PyTorch library.

TABLE VI  
THE MEAN DICE SCORES OF DIFFERENT MODELS FOR 30 IMAGES BETWEEN SALIENCY MAPS AND SEGMENTATION. THE SEGMENTATION IS OBTAINED BY REMOVING THE BACKGROUND FROM AN IMAGE.

Network	Dice Score
VGG19	0.609
SqueezeNet	0.586
<b>Proposed(1<sup>st</sup> run)</b>	<b>0.630</b>
Proposed(2 <sup>nd</sup> run)	0.629
Proposed(3 <sup>rd</sup> run)	0.627
Proposed(mean±std)	0.6287±0.0015

Saliency maps for four randomly selected ImageNet images from different models are shown in Figure 7. Visually, for all images, the saliency maps of the proposed model are sharper, and the brightest points more tightly conform to the objects, compared to VGG19 and SqueezeNet. In addition, we also quantify the quality of saliency maps. First, we threshold the saliency map by setting the  $q$ -percentile brightest pixels as one and the rest as zero to get a segmentation map. Then, we use the Dice score ( $\frac{2|X \cap Y|}{X \cup Y}$ ) [83] between the segmentation of an object and a saliency map to measure their similarity. This metric by and large can reflect the sharpness and accuracy of a saliency map. The higher the score is, the better interpretability a model has. The obtained segmentation and saliency maps are put in Part C of supplementary materials for conciseness. Figure 8 shows the Dice scores for four objects concerning different percentiles and models. The percentile range is from top-1% to top-20% with a step of 1%. We find that the proposed model achieves the highest Dice scores over common-newt, brambling, and quail images. For the corn image, the proposed model is comparable to VGG19 but much better than SqueezeNet.

Furthermore, we make a dataset comprising 30 images and their segmentation maps, by randomly selecting images from the ImageNet validation set and manually removing their background. For each pair, we record the maximum Dice score associated with a certain percentile. The mean Dice scores for 30 images are shown in Table VI. The detailed Dice scores for each image are summarized in Table I of Part D in supplementary materials. There are two highlights from Table VI. First, the Dice scores of our models surpass those of competitors with a considerable margin, which implies that the saliency maps generated by our model are of higher quality than those of competitors. The second highlight is that our model results are pretty consistent with one another, where the variance among models is only 0.0015. To highlight the improvements made by the proposed model, we conduct the paired t-test between the proposed model and the competitor. The null hypothesis is that the pairwise difference in Dice scores between two models has a mean equal to zero. The test decisions for all pairs are shown in Table II in Appendix C), where all decisions reject the null hypothesis at the default 5% significance level. This suggests that the improvement by our

model is significant.

## VI. DISCUSSION

In [84], it was demonstrated that the ResNet topology is also intrinsically the densely connected topology. Suppose  $X_l = H_l(R_{l-1})$ , which is the output of the  $l^{\text{th}}$  layer, and  $R_0 = X_0$ ,

$$\begin{aligned}
 X_l &= H_l(R_{l-1}) = H_l(H_{l-1}(R_{l-2}) + R_{l-2}) \\
 &= H_l(H_{l-1}(R_{l-2}) + H_{l-2}(R_{l-3}) + R_{l-3}) \\
 &= H_l\left(\sum_{i=0}^{l-1} H_i(R_{i-1}) + R_0\right) \\
 &= H_l\left(\sum_{i=0}^{l-1} X_i + X_0\right) \\
 &= H_l(X_0 + X_1 + \dots + X_{l-1}).
 \end{aligned} \tag{35}$$

Therefore, our theoretical results on the densely connected topology can be somehow extended to the ResNet topology.

In [29], ResNet is interpreted as an ensemble of many paths of different lengths, and an ablation study shows that deleting a single layer does not affect the performance significantly. In light of ensemble behavior, as shown in Figure 9, given the depth  $L$ , there are  $2^L$  implicit paths connecting the input and output in ResNet, while for the proposed network, the number of implicit paths is  $L + 1$ . Furthermore, in ResNet, every layer has an equal chance of being passed or not passed. However, implicit paths of the proposed topology pass earlier layers more than later layers. For example, in Figure 9(b), only one path connects  $H_2$ , but three paths connect  $H_0$ . We conduct the ablation study on the proposed network with  $k=180$  and depth = 41 from Table V. We set the outputs of the first and fifth layers of each block as zeros respectively and examine the performance of the network on the test set. Because we have obtained three models from repetitive experiments, the ablation is repeated three times. The results are shown in Table VII. We can see that undoing the first layer of each block has a significant impact, which causes only 31.84% accuracy. In contrast, the model with the fifth layer of each block being undone still has the classification accuracy of 61.91%.

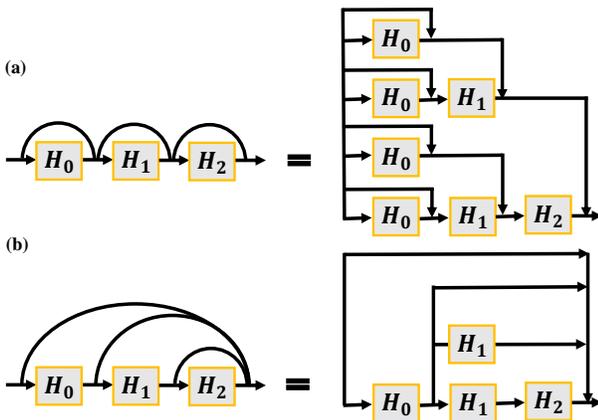


Fig. 9. (a) ResNet and its unraveled view; (b) the proposed topology and its unraveled view. The operation at the joints of black lines is summation.

TABLE VII

PERFORMANCE BY UNDOING DIFFERENT LAYERS IN THE PROPOSED MODEL TO MANIFEST THE RELATIVE IMPORTANCE OF EACH LAYER

	Original	Undo Layer 1	Undo Layer 5
Accuracy (%)	$77.1 \pm 0.06$	$31.84 \pm 1.81$	$61.91 \pm 5.52$

## VII. CONCLUSION

In this study, we have theoretically demonstrated the expressivity and generalizability of skip connections in deep learning, with an emphasis on the proposed topology. Then, we have performed comprehensive prediction and classification experiments to corroborate our theoretical findings that the networks of the proposed topology enjoy good expressivity and generalizability. Furthermore, we have also shown that the proposed model embraces improved interpretability in terms of saliency maps and layer importance. We have shared our code and prepared images in <https://github.com/FengleiFan/SparseShortcutTopology>. Future research directions can be put into exploring the utility of network equivalency in neural architecture search studies.

## REFERENCES

- [1] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [2] G. E. Dahl, D. Yu, L. Deng, and A. Acero, "Context-dependent pre-trained deep neural networks for large-vocabulary speech recognition," *IEEE Transactions on audio, speech, and language processing*, vol. 20, no. 1, pp. 30–42, 2011.
- [3] A. Kumar, O. Irsoy, P. Ondruska, M. Iyyer, J. Bradbury, I. Gulrajani, V. Zhong, R. Paulus, and R. Socher, "Ask me anything: Dynamic memory networks for natural language processing," in *International conference on machine learning*, pp. 1378–1387, PMLR, 2016.
- [4] H. Chen, Y. Zhang, M. K. Kalra, F. Lin, Y. Chen, P. Liao, J. Zhou, and G. Wang, "Low-dose ct with a residual encoder-decoder convolutional neural network," *IEEE transactions on medical imaging*, vol. 36, no. 12, pp. 2524–2535, 2017.
- [5] G. Wang, "A perspective on deep imaging," *Ieee Access*, vol. 4, pp. 8914–8924, 2016.
- [6] M. Anthimopoulos, S. Christodoulidis, L. Ebner, A. Christe, and S. Mougiakakou, "Lung pattern classification for interstitial lung diseases using a deep convolutional neural network," *IEEE transactions on medical imaging*, vol. 35, no. 5, pp. 1207–1216, 2016.
- [7] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in neural information processing systems*, vol. 25, pp. 1097–1105, 2012.
- [8] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2818–2826, 2016.
- [9] M. Lin, Q. Chen, and S. Yan, "Network in network," *International Conference on Learning Representations*, 2014.
- [10] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *International Conference on Learning Representations*, 2015.
- [11] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- [12] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4700–4708, 2017.
- [13] X. Chu, B. Zhang, H. Ma, R. Xu, and Q. Li, "Fast, accurate and lightweight super-resolution with neural architecture search," *arXiv preprint arXiv:1901.07261*, 2019.
- [14] B. Hariharan, P. Arbeláez, R. Girshick, and J. Malik, "Hypercolumns for object segmentation and fine-grained localization," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 447–456, 2015.

- [15] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 12, pp. 2481–2495, 2017.
- [16] R. K. Srivastava, K. Greff, and J. Schmidhuber, "Training very deep networks," in *Proceedings of the 28th International Conference on Neural Information Processing Systems-Volume 2*, pp. 2377–2385, 2015.
- [17] G. Larsson, M. Maire, and G. Shakhnarovich, "Fractalnet: Ultra-deep neural networks without residuals," *International Conference on Learning Representations*, 2017.
- [18] K.-I. Funahashi, "On the approximate realization of continuous mappings by neural networks," *Neural networks*, vol. 2, no. 3, pp. 183–192, 1989.
- [19] K. Hornik, M. Stinchcombe, and H. White, "Multilayer feedforward networks are universal approximators," *Neural networks*, vol. 2, no. 5, pp. 359–366, 1989.
- [20] L. Szymanski and B. McCane, "Deep networks are effective encoders of periodicity," *IEEE transactions on neural networks and learning systems*, vol. 25, no. 10, pp. 1816–1827, 2014.
- [21] D. Rolnick and M. Tegmark, "The power of deeper networks for expressing natural functions," *International Conference on Learning Representations*, 2018.
- [22] H. N. Mhaskar and T. Poggio, "Deep vs. shallow networks: An approximation theory perspective," *Analysis and Applications*, vol. 14, no. 06, pp. 829–848, 2016.
- [23] R. Eldan and O. Shamir, "The power of depth for feedforward neural networks," in *Conference on learning theory*, pp. 907–940, PMLR, 2016.
- [24] S. Liang and R. Srikant, "Why deep neural networks for function approximation?," in *International Conference on Learning Representations*, 2017.
- [25] V. Tikhomirov, "On the representation of continuous functions of several variables as superpositions of continuous functions of one variable and addition," in *Selected Works of AN Kolmogorov*, pp. 383–387, Springer, 1991.
- [26] Z. Lu, H. Pu, F. Wang, Z. Hu, and L. Wang, "The expressive power of neural networks: a view from the width," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pp. 6232–6240, 2017.
- [27] F. Fan, J. Xiong, and G. Wang, "Universal approximation with quadratic deep networks," *Neural Networks*, vol. 124, pp. 383–392, 2020.
- [28] H. Lin and S. Jegelka, "Resnet with one-neuron hidden layers is a universal approximator," *Advances in Neural Information Processing Systems*, vol. 31, pp. 6169–6178, 2018.
- [29] A. Veit, M. Wilber, and S. Belongie, "Residual networks behave like ensembles of relatively shallow networks," in *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pp. 550–558, 2016.
- [30] T. L. Liu, M. Chen, M. Zhou, S. Du, E. Zhou, and T. Zhao, "Towards understanding the importance of shortcut connections in residual networks," *Advances in neural information processing systems*, 2019.
- [31] F. He, T. Liu, and D. Tao, "Why resnet works? residuals generalize.," *IEEE transactions on neural networks and learning systems*, vol. 31, pp. 5349–5362, 2020.
- [32] E. Kang, H. J. Koo, D. H. Yang, J. B. Seo, and J. C. Ye, "Cycle-consistent adversarial denoising network for multiphase coronary ct angiography," *Medical physics*, vol. 46, no. 2, pp. 550–562, 2019.
- [33] C. You, G. Li, Y. Zhang, X. Zhang, H. Shan, M. Li, S. Ju, Z. Zhao, Z. Zhang, W. Cong, *et al.*, "Ct super-resolution gan constrained by the identical, residual, and cycle learning ensemble (gan-circle)," *IEEE transactions on medical imaging*, vol. 39, no. 1, pp. 188–203, 2019.
- [34] B. Hamann and J.-L. Chen, "Data point selection for piecewise linear curve approximation," *Computer Aided Geometric Design*, vol. 11, no. 3, pp. 289–301, 1994.
- [35] L. Zhu, R. Deng, M. Maire, Z. Deng, G. Mori, and P. Tan, "Sparsely aggregated convolutional networks," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 186–201, 2018.
- [36] J. Schmidt-Hieber, "The kolmogorov–arnold representation theorem revisited," *Neural Networks*, vol. 137, pp. 119–126, 2021.
- [37] P. L. Bartlett, D. J. Foster, and M. Telgarsky, "Spectrally-normalized margin bounds for neural networks," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pp. 6241–6250, 2017.
- [38] B. Neyshabur, R. Tomioka, and N. Srebro, "Norm-based capacity control in neural networks," in *Conference on Learning Theory*, pp. 1376–1401, PMLR, 2015.
- [39] B. Neyshabur, S. Bhojanapalli, and N. Srebro, "A pac-bayesian approach to spectrally-normalized margin bounds for neural networks," in *International Conference on Learning Representations*, 2018.
- [40] L. Wu, C. Ma, and W. E, "How sgd selects the global minima in over-parameterized learning: A dynamical stability perspective," in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pp. 8289–8298, 2018.
- [41] Z. Allen-Zhu, Y. Li, and Y. Liang, "Learning and generalization in overparameterized neural networks, going beyond two layers," *arXiv preprint arXiv:1811.04918*, 2018.
- [42] A. Brutzkus, A. Globerson, E. Malach, and S. Shalev-Shwartz, "Sgd learns over-parameterized networks that provably generalize on linearly separable data," in *International Conference on Learning Representations*, 2018.
- [43] B. Neyshabur, Z. Li, S. Bhojanapalli, Y. LeCun, and N. Srebro, "The role of over-parametrization in generalization of neural networks," in *International Conference on Learning Representations*, 2018.
- [44] P. Nakkiran, G. Kaplun, Y. Bansal, T. Yang, B. Barak, and I. Sutskever, "Deep double descent: Where bigger models and more data hurt," in *International Conference on Learning Representations*, 2019.
- [45] V. Vovk, "Kernel ridge regression," in *Empirical inference*, pp. 105–116, Springer, 2013.
- [46] A. Jacot, F. Gabriel, and C. Hongler, "Neural tangent kernel: Convergence and generalization in neural networks," in *Advances in neural information processing systems*, pp. 8571–8580, 2018.
- [47] S. Arora, S. S. Du, W. Hu, Z. Li, R. R. Salakhutdinov, and R. Wang, "On exact computation with an infinitely wide neural net," in *Advances in Neural Information Processing Systems*, pp. 8141–8150, 2019.
- [48] D. Harrison Jr and D. L. Rubinfeld, "Hedonic housing prices and the demand for clean air," *Journal of environmental economics and management*, vol. 5, no. 1, pp. 81–102, 1978.
- [49] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International conference on machine learning*, pp. 448–456, PMLR, 2015.
- [50] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *The journal of machine learning research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [51] G. Huang, Y. Sun, Z. Liu, D. Sedra, and K. Q. Weinberger, "Deep networks with stochastic depth," in *European conference on computer vision*, pp. 646–661, Springer, 2016.
- [52] Z. Huang, S. Liang, M. Liang, and H. Yang, "Dianet: Dense-and-implicit attention network.," in *AAAI*, pp. 4206–4214, 2020.
- [53] H. Kabir, M. Abdar, S. M. J. Jalali, A. Khosravi, A. F. Atiya, S. Nahavandi, and D. Srinivasan, "Spinalnet: Deep neural network with gradual input," *arXiv preprint arXiv:2007.03347*, 2020.
- [54] G. Franchi, A. Bursuc, E. Aldea, S. Dubuisson, and I. Bloch, "Encoding the latent posterior of bayesian neural networks for uncertainty quantification," *arXiv preprint arXiv:2012.02818*, 2020.
- [55] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4510–4520, 2018.
- [56] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *International Conference on Machine Learning*, pp. 6105–6114, PMLR, 2019.
- [57] Y. Chen, H. Fan, B. Xu, Z. Yan, Y. Kalantidis, M. Rohrbach, S. Yan, and J. Feng, "Drop an octave: Reducing spatial redundancy in convolutional neural networks with octave convolution," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 3435–3444, 2019.
- [58] I. Bello, "Lambdanetworks: Modeling long-range interactions without attention," in *International Conference on Learning Representations*, 2020.
- [59] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7132–7141, 2018.
- [60] Y. Li, Z. Kuang, Y. Chen, and W. Zhang, "Data-driven neuron allocation for scale aggregation networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 11526–11534, 2019.
- [61] K. Han, Y. Wang, Q. Tian, J. Guo, C. Xu, and C. Xu, "Ghostnet: More features from cheap operations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1580–1589, 2020.
- [62] S. Xie, A. Kirillov, R. Girshick, and K. He, "Exploring randomly wired neural networks for image recognition," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1284–1293, 2019.

- [63] I. Bello, "Lambdanetworks: Modeling long-range interactions without attention," in *Submitted to International Conference on Learning Representations*, 2021. under review.
- [64] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255, Ieee, 2009.
- [65] X. Zhang, X. Zhou, M. Lin, and J. Sun, "Shufflenet: An extremely efficient convolutional neural network for mobile devices," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 6848–6856, 2018.
- [66] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8697–8710, 2018.
- [67] E. Real, A. Aggarwal, Y. Huang, and Q. V. Le, "Regularized evolution for image classifier architecture search," in *Proceedings of the aaai conference on artificial intelligence*, vol. 33, pp. 4780–4789, 2019.
- [68] C. Liu, B. Zoph, M. Neumann, J. Shlens, W. Hua, L.-J. Li, L. Fei-Fei, A. Yuille, J. Huang, and K. Murphy, "Progressive neural architecture search," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 19–34, 2018.
- [69] H. Liu, K. Simonyan, and Y. Yang, "Darts: Differentiable architecture search," in *International Conference on Learning Representations*, 2018.
- [70] B. Wu, X. Dai, P. Zhang, Y. Wang, F. Sun, Y. Wu, Y. Tian, P. Vajda, Y. Jia, and K. Keutzer, "Fbnet: Hardware-aware efficient convnet design via differentiable neural architecture search," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 10734–10742, 2019.
- [71] I. Radosavovic, R. P. Kosaraju, R. Girshick, K. He, and P. Dollár, "Designing network design spaces," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10428–10436, 2020.
- [72] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou, "Training data-efficient image transformers & distillation through attention," in *International Conference on Machine Learning*, pp. 10347–10357, PMLR, 2021.
- [73] G. Wang, K. Wang, and L. Lin, "Adaptively connected neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1781–1790, 2019.
- [74] J. Fang, Y. Sun, Q. Zhang, Y. Li, W. Liu, and X. Wang, "Densely connected search space for more flexible neural architecture search," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10628–10637, 2020.
- [75] Q. Wang, B. Wu, P. Zhu, P. Li, W. Zuo, and Q. Hu, "Eca-net: Efficient channel attention for deep convolutional neural networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11534–11542, 2020.
- [76] F. Fan, J. Xiong, M. Li, and G. Wang, "On interpretability of artificial neural networks: A survey," *arXiv preprint arXiv:2001.02522*, 2020.
- [77] B.-J. Hou and Z.-H. Zhou, "Learning with interpretable structure from gated rnn," *IEEE transactions on neural networks and learning systems*, vol. 31, no. 7, pp. 2267–2279, 2020.
- [78] A. B. Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. García, S. Gil-López, D. Molina, R. Benjamins, *et al.*, "Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai," *Information Fusion*, vol. 58, pp. 82–115, 2020.
- [79] S. Srinivas and F. Fleuret, "Full-gradient representation for neural network visualization," in *Advances in Neural Information Processing Systems*, pp. 4124–4133, 2019.
- [80] D. Smilkov, N. Thorat, B. Kim, F. Viégas, and M. Wattenberg, "Smoothgrad: removing noise by adding noise," *International conference on machine learning*, 2017.
- [81] M. Sundararajan, A. Taly, and Q. Yan, "Axiomatic attribution for deep networks," *International conference on machine learning*, 2017.
- [82] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, "Squeezenet: Alexnet-level accuracy with 50x fewer parameters and < 0.5 mb model size," *International Conference on Learning Representations*, 2017.
- [83] L. R. Dice, "Measures of the amount of ecologic association between species," *Ecology*, vol. 26, no. 3, pp. 297–302, 1945.
- [84] W. Wang, X. Li, T. Lu, and J. Yang, "Mixed link networks," in *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pp. 2819–2825, 2018.