Automatic Sparse Connectivity Learning for Neural Networks

Zhimin Tang, Linkai Luo*, Bike Xie, Yiyu Zhu, Rujie Zhao, Lvqing Bi, Chao Lu*

Abstract—Since sparse neural networks usually contain many zero weights, these unnecessary network connections can potentially be eliminated without degrading network performance. Therefore, well-designed sparse neural networks have the potential to significantly reduce FLOPs and computational resources. In this work, we propose a new automatic pruning method -Sparse Connectivity Learning (SCL). Specifically, a weight is reparameterized as an element-wise multiplication of a trainable weight variable and a binary mask. Thus, network connectivity is fully described by the binary mask, which is modulated by a unit step function. We theoretically prove the fundamental principle of using a straight-through estimator (STE) for network pruning. This principle is that the proxy gradients of STE should be positive, ensuring that mask variables converge at their minima. After finding Leaky ReLU, Softplus, and Identity STEs can satisfy this principle, we propose to adopt Identity STE in SCL for discrete mask relaxation. We find that mask gradients of different features are very unbalanced, hence, we propose to normalize mask gradients of each feature to optimize mask variable training. In order to automatically train sparse masks, we include the total number of network connections as a regularization term in our objective function. As SCL does not require pruning criteria or hyper-parameters defined by designers for network layers, the network is explored in a larger hypothesis space to achieve optimized sparse connectivity for the best performance. SCL overcomes the limitations of existing automatic pruning methods. Experimental results demonstrate that SCL can automatically learn and select important network connections for various baseline network structures. Deep learning models trained by SCL outperform the state-of-theart human-designed and automatic pruning methods in sparsity, accuracy, and FLOPs reduction.

Index Terms—neural networks, model compression, model pruning, sparse connectivity learning, trainable binary mask.

I. INTRODUCTION

D ESPITE the great success in improving neural networks and learning systems [1]–[6], state-of-the-art deep neural networks usually consist of dozens of stacked layers and a huge number of parameters. It is difficult to deploy these over-parameterized and over-redundant neural networks on resource-constrained computing platforms [7]. To address this

*Corresponding authors: L. Luo and C. Lu

challenge, network pruning has received great attention. Network pruning tends to remove unimportant trainable parameters in a neural network architecture while maintaining its high accuracy. Effective network pruning leads to less computing operations, memory usage, and power consumption with little performance degeneration. After performing network pruning, sparse models are often implemented in hardware (*e.g.*, GPUs, ASICs¹ or FPGAs²) for AI acceleration.

1

Network pruning can be divided into two categories: humandesigned pruning and automatic pruning. The former requires some pruning criteria or hyper-parameters defined by designers (e.g., importance measure or pruning threshold), while automatic pruning generates optimized sparse networks with little human intervention. Human-designed pruning usually consists of three steps: (1) training weight parameters in a selected baseline neural network, (2) eliminating unimportant network connections based on designer-defined criteria or hyper-parameters, and (3) training weight parameters again in this pruned network architecture. Human-designed network pruning can be further divided into two types: unstructured [7]–[13] or structured [14]–[26]. It has been reported that performing unstructured pruning on deep neural networks does not cause much loss of accuracy. On the other hand, structured pruning especially filter [14]–[19] and channel [20]–[24] pruning, has been used to accelerate neural networks in general hardware platforms (i.e., GPUs). Since zero-masked feature maps can be deleted, less computational cost is required after network pruning. Note that some existing neural networks with multi-branch or multi-group structures may be considered as human-designed sparse architectures (e.g., Inception [1], ShuffleNet [27], MobileNet [28], ResNeXt [29]), even though these sparse architectures do not involve network pruning.

Up to date, three major bottlenecks have hindered the use of human-designed pruning methods to generate sparse networks. First, although human-designed pruning can provide good network performance under a low pruning rate, the network performance under a high pruning rate is severely degraded. Second, there is a lack of appropriate methods to effectively prune network connections for high-compression and highperformance neural networks. In most human-designed pruning methods, network connections are pruned based on the assumption of "smaller-norm-less-important". Network connections with smaller weights are generally considered trivial and eliminated during network pruning. However, researchers have found that sometimes the amount of information in smaller

Z. Tang and L. Luo are with Department of Automation, Xiamen University, Xiamen, 361102, China. Tang is a visiting scholar at the ECE Department of SIUC from 2018 to 2020. (e-mail: luolk@xmu.edu.cn)

B. Xie and Y. Zhu are with Kneron Inc. San Diego, CA, 92121, USA.

L. Bi is with School of Physics and Telecommunication Engineering, Research Center for Intelligent Information and Communication Technology, Yulin Normal University, Yulin 537000, Guangxi, China.

R. Zhao and C. Lu are with the Department of Electrical and Computer Engineering, Southern Illinois University Carbondale (SIUC), Carbondale, IL, 62091, USA. (e-mail: chaolu@siu.edu)

¹ https://www.kneron.com/solutions/soc/

² https://www.xilinx.com/applications/megatrends/machine-learning.html

weights is important and cannot be ignored [30]. Therefore, considering the variability of loss function sensitivity with respect to different weights, eliminating network connections with smaller weights does not guarantee a slight decrease in network performance. Third, the biggest weakness is the need for designer-defined pruning criteria (e.g., L_1 [7] or L_2 [16] norm of filters) and hyper-parameters (e.g., pruning threshold or ratio) for each network layer during network pruning. Because the selection of pruning criteria or hyperparameters is manually determined, it heavily depends on the designer's prior experience and varies from application to application. Consequently, since the pruning criteria and hyper-parameters are not guaranteed to be the best choice, the resulting network connectivity and performance are not optimal. Note that even though weight importance estimation methods have been proposed in [31], [32], layer-wise hyperparameters are still needed to determine a pruning threshold or ratio for all layers during network pruning. Unfortunately, as the optimal pruning threshold or ratio varies with local network structures, the use of only one hyper-parameter for all layers leads to inferior pruning performance. Recently, it is proposed that a proper criterion may be selected from a set of designerdefined criteria to learn pruning criteria for each network layer [33]. Yet, the pruning performance of this method still depends on the quality of candidate criteria defined by designers.

To get rid of shortcomings of human-designed pruning, researchers have investigated automatic pruning, which trains sparse network connectivity through task-aware loss or a sparse regularization term [34]-[38]. Note automatic pruning does not require designer-defined pruning criteria or layer-wise hyper-parameters. Louizos et al. [34] train sparse neural networks through L_0 -norm regularization. They use the Gumbel-Softmax trick [39], [40] (also known as a concrete distribution) and apply gates to weights for connectivity training. Note a stochastic gate produces zero connectivity only when the probability is zero, which is almost impossible to reach during network training. To address this problem, threshold operation is applied on gates to made zero connectivity. The drawback of [34] is that the expected L_0 -norm of stochastic training does not reflect the L_0 -norm of deterministic inference. Therefore, as will be discussed in Section V-D3, although the expected L_0 norm of weights is significantly reduced, it is still not low enough to produce sparse connections. Kang et al. [35] propose soft channel pruning (SCP), which assumes that feature maps follow a Gaussian distribution and features are pruned if the cumulative density function is larger than a certain threshold. The Gumbel-Softmax trick is used to tackle the nondifferentiable Bernoulli distribution sampling. Unfortunately, the assumption of Gaussian distribution for feature maps is too strict to derive accurate gradients. Moreover, SCP shows good pruning performance in network layers that are followed by both batch normalization (BN) and ReLU. According to evaluation results in [35], its pruning performance is severely degraded if only BN exists. Since many neural networks do not include both BN and ReLU, the application scope of SCP is limited. Herrmann et al. [36] jointly consider conditional computation and network pruning. The Gumbel-Softmax trick is used to relax the discrete masks to a continuous form. The researchers focus on dynamic inference using conditional computation. Depending on network inputs, masks are dynamically applied on channels. Unfortunately, these masks are data-dependent rather than being fixed, and the dynamically pruned structures are not friendly to hardware implementation. Huang et al. [37] introduce a series of non-negative scaling factors that are associated with neural network connectivity. To encourage sparse network connectivity, these scaling factors are penalized by an L_1 -norm regularization term. During the training process, stochastic gradient descent (SGD) and a proposed stochastic Accelerated Proximal Gradient (APG) are used to train weight parameters and scaling factor parameters, respectively. The scaling factor parameters are converted into scaling factors through a soft-threshold operation. Xiao et al. [38] utilize STE to relax binary masks. Even though empirical experiments have shown potential, the fundamental principle of using STE in network pruning has not been explored.

From the above discussion, it is clear that existing network pruning methods, either human-designed or automatic, do not fully address the requirement of highly effective spare connectivity learning. It is attractive to develop new pruning methods to overcome the drawbacks and limitations of existing pruning methods. In this work, we propose a new automatic network pruning technique - sparse connectivity learning (SCL). This work makes the following contributions:

- We theoretically prove the fundamental principle of using STE for network pruning. After finding Leaky ReLU, Softplus, and Identity STEs can satisfy this principle, we propose to adopt Identity STE in SCL for discrete mask relaxation. Thus, SCL guarantees the convergence of mask variables at their minima.
- We observe that mask gradients on different features have a wide range of magnitudes and hence are unbalanced. Therefore, we propose to normalize mask gradients of each feature to optimize mask variable training.
- The pruning principle of our proposed SCL method is the significance of weight, instead of the magnitude of weight. SCL can automatically learn and determine critical network connections of baseline networks (*e.g.*, DenseNets, ResNets, VGGs, EfficientNets, and RNNs).
- SCL enables highly effective weight-level sparsity learning on neural networks under high pruning rates. Experimental results in the MNIST, CIFAR-10, CIFAR-100, ImageNet, and WikiText-2 datasets demonstrate the enhanced performance of SCL-induced sparse neural networks than the state-of-the-art network pruning methods in the literature.
- SCL automatically learns optimized neural network connectivity in a task-aware manner, ensuring that performance-sensitive network connections are ultimately preserved. As it does not need any designer-defined pruning criteria or layer-wise pruning hyper-parameters, SCL gets rid of the limitation of human designers. Compared to the state-of-the-art pruning methods, experimental results demonstrate that SCL results in high-performance neural networks with higher sparsity and fewer FLOPs.

Since this work focuses on the algorithm aspect of net-

work pruning, the hardware implementation of spare models generated from our proposed pruning algorithm is beyond the scope of this paper. Sparse neural network models generated from SCL can be implemented in FPGAs or Kneron edge AI hardware products. This paper is organized as follows. Section II introduces related works about pruning criteria and binary mask relaxation. Section III describes the proposed automatic sparse connectivity learning theory. Section IV introduces the baseline network architectures and experimental setups. Section V demonstrates various experimental results and comparisons with the state-of-the-art works in the literature. Section VI concludes the paper.

II. RELATED WORK

A. Human-designed pruning criteria

The performance of human-designed pruning methods depends to a large extent on the quality of pruning criteria defined by designers. Unimportant network connections are removed based on the importance measure, which usually follows the assumption that smaller norms are less important. The criterion for unstructured pruning is the absolute value of weights [7], [8]. The pruning criterion for structured pruning is the *L*1-norm [14] or *L*2-norm [16] of filter or channel weights. The magnitude of scaling factors in batch normalization is also proposed as the channel pruning criterion [20]. He *et al.* [19] do not agree that smaller filters are less important, but propose that the contribution of median filters is relatively small, because they can be represented by other filters.

B. Gumbel-Softmax trick

Automatic pruning involves training binary masks. To address the non-differentiation problem, the Gumbel-Softmax trick [39], [40] has been used to relax binary masks [34]– [36]. This trick uses gradient methods to train discrete random variables. In this trick, discrete values produced by argmax are encoded in a one-hot vector, and the random sampling f(U)is expressed as

$$f(U) = argmax(\log \alpha + G(U)) \tag{1}$$

$$G(U) = -\log(-\log U), U \sim \mathcal{U}(0, 1)$$
(2)

where the argmax function finds the argument corresponding to the maximum value, α is a set of unnormalized parameters α_k and the probability of outcome k is $\alpha_k / \sum_i \alpha_i$. Each element in the vector G(U) obeys a Gumbel distribution. As a binary mask has two possible discrete values, mask re-parameterization is regarded as a special case. Due to the sparsity consideration, there is a high possibility that mask sampling results should be trained to be zero. As a result, sparse binary masks are obtained. However, the discrete argmax operation can not be trained by gradient methods. Therefore, one way to circumvent this is to relax the argmax operation by replacing it with softmax [39] and [40], as

$$f(U) = softmax((\log \alpha + G(U))/\tau)$$
(3)

where τ is a hyper-parameter of the softmax function to control relaxation. When $\tau \to 0$, the softmax function becomes the argmax function. Thus, Eq. (3) provides a differentiable form, which is able to train and optimize by gradient methods.

The Gumbel-Softmax trick has two limitations. First, its gradient estimations are biased with respect to the gradients of discrete connectivity. According to Eq. (3), this trick leads to unbiased gradients only when the hyper-parameter τ tends to 0. However, because the value of τ is used to balance bias and variance in practice, τ is rarely chosen to be close to 0. As a result, the studies in [34]–[36] use biased gradients, which may not meet the essential condition for convergence (*i.e.*, low-biased gradient). Second, the training process is stochastic, but the inference is deterministic. Therefore, the expected number of connections of stochastic training does not reflect the number of connections of deterministic inference. As shown in Section V-D3, the reduction of L_0 -norm during training does not mean sparser network connections.

C. Straight-through estimator (STE) trick

Due to the derivatives of the binarization function are mostly zero, training variables can not update using gradientbased optimization methods. STE is a trick of using proxy gradients in back-propagation [41]. In this trick, zero gradients of a discrete function are replaced by the derivative of a (sub)differentiable function. Derivatives of ReLU STE and Clipped ReLU STE were used for network quantization [42]-[44]. Yin et al. [45] referred to proxy gradients as coarse gradients, and studied the STE properties for network quantization. [45] assumes that network inputs obey a normal distribution, yet, this assumption has not been validated for network pruning. Later, based on [45], Xiao et al. [38] proposed to use Leaky ReLU STE and Softplus STE to relax binary masks. Instead of performing network pruning, Hinton et al. [46] proposed to use Identity STE for binary neuron training. The (sub)differentiable functions of five existing STEs are expressed as

$$\sigma_{relu}(x) = \max(0, x) \tag{4}$$

$$\sigma_{clipped_relu}(x) = \min(\max(0, x), \alpha), \alpha > 0$$
(5)

$$\sigma_{leaky_relu}(x) = \max(x, \alpha \cdot x), 0 < \alpha < 1$$
(6)

$$\sigma_{softplus}(x) = \log(1 + \exp(x)) \tag{7}$$

$$\sigma_{identity}(x) = x \tag{8}$$

The STE trick for network pruning has two limitations. First, as Xiao *et al.* [38] treat connectivity as a hyperparameter, bi-level optimization is used to update weights and masks separately. Because it involves two-pass calculations to complete an update, the pruning optimization method is too cumbersome and computationally expensive. Second, although the use of STE for network pruning has achieved empirical success in [38], yet, it is based on the assumption that network inputs obey a normal distribution. Unfortunately, so far, the fundamental principle of using STE for network pruning is still unclear, which is the focus of this work.



Fig. 1. Automatic sparse connectivity learning. Weight re-parameterization for sparse convolution, Identity STE for masked weight training, Identity STE for binary mask relaxation, and mask gradient normalization. The elliptical symbols, 'Weight_V' and 'Mask_V', refer to trainable variables, the rectangular symbols indicate intermediate tensors, and the circles indicate calculation operators. Dotted lines mean gradient process during backward propagation.

III. AUTOMATIC SPARSE CONNECTIVITY LEARNING

The proposed automatic sparse connectivity learning is briefly described below and illustrated in Figure 1. First, we represent the weight as a multiplication of a weight variable and a binary mask (0 or 1). 1 indicates a network connection, while 0 indicates no network connection. Thus, the connectivity of a neural network can be fully described by the binary mask, which will be further modulated by a unit step function on a mask variable. Second, during the training process, we will decay the network connectivity to gradually push the elements of the binary mask towards zero (*i.e.*, no network connection). Finally, without the need of designer-defined criteria or layer-wise pruning hyper-parameters, unimportant network connections will be automatically discovered and removed. The design details of each step will be elaborated in the following subsections.

A. Weight re-parameterization

We design to learn a binary mask that is a standard convolution layer and is related to network connectivity. A fully connected layer is a special case of an input with a feature size of 1×1 convolved with a kernel with a size of 1×1 . The base cell of a RNN is actually composed of several fully connected modules. We model this convolution operation as

$$\mathbf{y} = \mathbf{w} * \mathbf{x} \tag{9}$$

where the convolutional kernel w is weight, x is input, and y is output. * annotates the convolution operation. As shown in Figure 1, a weight is re-parameterized by a weight variable \tilde{w} and a binary mask m. Weight re-parameterization is formulated as

$$\mathbf{w} = \widetilde{\mathbf{w}} \odot \mathbf{m} \tag{10}$$

where \odot is the element-wise multiplication. Mask is binarized from the trainable mask variable \widetilde{m} using a unit step function. The binarization is formulated as

$$\mathbf{m} = H(\widetilde{\mathbf{m}}) = \begin{cases} 0, \widetilde{m} \le 0\\ 1, \widetilde{m} > 0 \end{cases}$$
(11)

Therefore, the elements of $\widetilde{\mathbf{w}}$ will be zero-masked if the corresponding elements in $\widetilde{\mathbf{m}}$ is non-positive.

B. Gradient redefinition for weight variables

According to the derivative chain rule and Eq. (10), the gradient of loss function \mathcal{L} with respect to $\widetilde{\mathbf{w}}$ is written as

$$\frac{\partial \mathcal{L}}{\partial \widetilde{\mathbf{w}}} = \frac{\partial \mathcal{L}}{\partial \mathbf{w}} \odot \mathbf{m}$$
(12)

As **m** is a sparse tensor, many elements of $\partial \mathcal{L}/\partial \widetilde{\mathbf{w}}$ are zero, indicating that these masked variables are not updated. Even though network connections that contribute less to precision can be ignored in the current global structure, they may not be negligible in future global structures. As a result, temporary zero masks do not mean unimportant, and it is worth keeping training for these zero-masked weight variables. In this work, the gradient of loss function \mathcal{L} with respect to a weight variable $\widetilde{\mathbf{w}}$ is redefined as

$$\frac{\partial \mathcal{L}}{\partial \widetilde{\mathbf{w}}} := \frac{\partial \mathcal{L}}{\partial \mathbf{w}}$$
(13)

where $\partial \mathcal{L}/\partial \mathbf{w}$ is obtained through back-propagation. Similar to [8], [16], even if some weight variables are temporarily zero-masked during training, they will update later.

C. Gradient redefinition for mask variables

Due to the derivatives of a unit step function are mostly zero, mask variables can not update using gradient-based optimization methods. To solve this problem, we investigate the fundamental principle of using STE for network pruning. Let us analyze a realizable case that uses the following loss function to update mask variables $\tilde{\mathbf{m}}$.

$$\min_{\tilde{\mathbf{m}}} \ell(\tilde{\mathbf{m}}) = \frac{1}{2} (H(\tilde{\mathbf{m}}) - \mathbf{m}^*)^2$$
(14)

where $\mathbf{H}(\tilde{\mathbf{m}})$ and \mathbf{m}^* represent the binarized mask values and optimal mask values, respectively. Note that the optimal solution of mask variables $\tilde{\mathbf{m}}$ is a region rather than a point. The gradient of $\tilde{\mathbf{m}}$ is expressed as

$$\frac{\partial \ell}{\partial \tilde{\mathbf{m}}} = \frac{\partial H(\tilde{\mathbf{m}})}{\partial \tilde{\mathbf{m}}} (H(\tilde{\mathbf{m}}) - \mathbf{m}^*)$$
(15)



Fig. 2. Proxy gradients of five existing STEs used for binarization relaxation. The X-axis and Y-axis represent mask variables $\tilde{\mathbf{m}}$ and proxy gradients $\partial H(\tilde{\mathbf{m}})/\partial \tilde{\mathbf{m}}$, respectively. The dead zones for ReLU and Clipped ReLU STEs are marked.

As $H(\tilde{m}_i)$ and m_i^* are binary values, there are three possible values (*i.e.*, 0, -1, and +1) for their difference $H(\tilde{m}_i) - m_i^*$. Hence, the gradient of \tilde{m}_i is expressed as

$$\frac{\partial \ell}{\partial \tilde{m}_{i}} = \begin{cases} \frac{\partial H(\tilde{m}_{i})}{\partial \tilde{m}_{i}} \cdot 0 , & if \quad H(\tilde{m}_{i}) = m_{i}^{*} \\ \frac{\partial H(\tilde{m}_{i})}{\partial \tilde{m}_{i}} \cdot (-1), & if \quad \tilde{m}_{i} \leq 0 \text{ and } m_{i}^{*} = 1 \\ \frac{\partial H(\tilde{m}_{i})}{\partial \tilde{m}_{i}} \cdot (+1), & if \quad \tilde{m}_{i} > 0 \text{ and } m_{i}^{*} = 0 \end{cases}$$
(16)

When STE is used to relax the binarization, $\partial H(\tilde{m}_i)/\partial \tilde{m}_i$ is replaced by proxy gradients. Correct proxy gradients should move \tilde{m}_i towards their optimal values during network pruning. The gradient should be zero when the optimal value is reached, indicating no further update. Based on the mechanism of gradient descent optimizers (*i.e.*, a variable is adjusted in the opposite direction of its gradient), \tilde{m}_i should move towards the negative direction of proxy gradients. Consequently, when STE is used for network pruning, positive values of $\partial H(\tilde{m}_i)/\partial \tilde{m}_i$ in Eq. (16) ensure mask variables converge at their minima. This fundamental principle is derived without any assumptions, and is just based on the mechanism of gradient descent optimizers. In contrast, a normal distribution is assumed for inputs in [38].

Figure 2 plots proxy gradients of the existing five STEs. We can see that three STEs (*i.e.*, Leaky ReLU, Softplus, and identity) can satisfy the fundamental principle of positive proxy gradients. In this work, the Identity STE is selected for simplicity. Since its proxy gradient is always 1 as shown in Figure 2(e), Eq. (16) is expressed as

$$\frac{\partial \ell}{\partial \tilde{m}_{i}} = \begin{cases} 0, & if \ H(\tilde{m}_{i}) = m_{i}^{*} \\ -1, & if \ \tilde{m}_{i} \le 0 \ and \ m_{i}^{*} = 1 \\ +1, & if \ \tilde{m}_{i} > 0 \ and \ m_{i}^{*} = 0 \end{cases}$$
(17)

As indicated in Eq. (17), when \tilde{m}_i reaches its optimal value, the gradient is zero. As a result, \tilde{m}_i does not update further. When $\tilde{m}_i \leq 0$ and $m_i^* = 1$, the gradient of -1 pushes \tilde{m}_i to be more positive. When $\tilde{m}_i > 0$ and $m_i^* = 0$, the gradient of +1 pushes \tilde{m}_i to be more negative. Thus, mask variables and network connectivity (from 1 to 0 or vice versa) can update during training. The Identity STE ensures that mask variables move towards and finally stabilize at their optimal values. According to Eq. (11) and Figure 2(e), we obtain

$$\frac{\partial \mathbf{m}}{\partial \widetilde{\mathbf{m}}} = \frac{\partial H(\widetilde{\mathbf{m}})}{\partial \widetilde{\mathbf{m}}} = 1 \tag{18}$$

which indicates that the differential of \tilde{m} is redefined as that of m. Thus, the gradient of loss function \mathcal{L} with respect to mask variables $\tilde{\mathbf{m}}$ is redefined as

$$\frac{\partial \mathcal{L}}{\partial \widetilde{\mathbf{m}}} = \frac{\partial \mathcal{L}}{\partial \mathbf{m}} \frac{\partial \mathbf{m}}{\partial \widetilde{\mathbf{m}}} := \frac{\partial \mathcal{L}}{\partial \mathbf{m}} = \frac{\partial \mathcal{L}}{\partial \mathbf{w}} \odot \widetilde{\mathbf{w}}$$
(19)

D. Mask gradient normalization

We observe a wide range of mask gradient magnitudes in various layers and channels of common neural networks. To visualize this observation, a pre-trained ResNet-110 is used as an example to plot mask gradient magnitudes in Figure 3. In Figure 3(a), layer-wise results show that the average mask gradient magnitude of the first and last layers (*e.g.*, 10^{-5}) is much higher than other layers (*e.g.*, 10^{-7}). Then, channel-wise results of two typical layers in Figures 3(b) and 3(c) show that mask gradient magnitudes fluctuate greatly between channels.

The influence of a wide range of mask gradient magnitudes is analyzed below. Let us review the training process of mask variables during network pruning. Through binarization (i.e., Eq. (11)), a positive mask variable means a mask state of 1, indicating that there is a network connection. If this network connection is pruned, the binary mask state should become 0. According to Eq. (11), the final value of this mask variable should be zero or negative. Assume two mask variables are initialized to the same positive value and eventually become zero. That is, two initially established network connections are eliminated after pruning. Under the same learning rate for all layers, the mask variable with a smaller gradient magnitude requires more training iterations to update until convergence. Therefore, it is difficult for mask variables with a wide range of mask gradient magnitudes to converge to their optimal solution, thereby degrading network pruning performance. To mitigate this problem, we propose to normalize mask gradients on different features. Mask gradients of each feature are normalized to the unit variance in each mini-batch. As shown in Eq. (20), mask gradients obtained through back-propagation are divided by a gradient scale s.

$$Norm(\frac{\partial \mathcal{L}}{\partial \mathbf{m}_{j}}) = \frac{\partial \mathcal{L}}{\partial \mathbf{w}_{j}} \odot \widetilde{\mathbf{w}}_{j} / (s+\epsilon) , \quad where$$

$$s = \sqrt{\sum_{\mathbf{x}_{b} \in \mathcal{B}} \sum_{w_{k} \in \mathbf{w}_{j}} \left(\frac{\partial \mathcal{L}(\mathbf{x}_{b})}{\partial w_{k}} \odot \widetilde{w}_{k}\right)^{2} / |\mathcal{B}| \cdot |\mathbf{w}_{j}|}$$
(20)

where \mathcal{B} and $|\mathcal{B}|$ denote the sampled mini-batch data and batch size, respectively. \mathbf{w}_j and $|\mathbf{w}_j|$ represent the weight and weight



Fig. 3. Mask gradient magnitudes of pre-trained ResNet-110. (b) and (c) show mask gradient magnitudes of channels in layer-1 and layer-105, respectively.

number of the *j*-th feature, respectively. ϵ is a small positive constant to avoid division by zero. Since it is not a zero-mean normalization, the sign of each mask gradient does not change, so the analysis derived in the previous subsection is still valid. Because the same gradients are produced on all masks by the connectivity decay introduced in the next subsection, the gradient normalization excludes mask gradients caused by the connectivity decay (*i.e.*, normalization in Eq. (20) processes the gradients produced by the first and third terms of Eq. (23)).

E. Connectivity decay for sparse connectivity learning

Through the proposed weight re-parameterization, the information of network connectivity is completely represented by the elements in the binary mask **m**. The degree of network connectivity is equal to the sum of all element values in **m**. We incorporate the degree of connectivity into an objective function to optimize network connectivity. The training objective function is expressed as

$$\mathcal{L} = \mathcal{C} + \lambda_1 \cdot \sum_{l=0}^{L} \sum_{i=0}^{|\mathbf{m}^{(l)}|} m_i^{(l)}$$
(21)

where C is the criteria to measure performance loss, $\mathbf{m}^{(l)}$ is the mask of *l*-th layer, *i* is the element index, $|\mathbf{m}^{(l)}|$ is the total number of elements in l-th layer, and L is the total number of layers. λ_1 is a hyper-parameter for connectivity decay. As the second term includes the information of network connectivity, so during training this term attenuation indicates connectivity decay. Connectivity decay can be viewed as a way of L_0 regularization of weight. There are two gradients that affect the training of mask variables. One gradient is due to the connectivity decay, which pushes mask variables moving towards the negative infinity direction all the time. The other gradient is due to the performance loss minimization, which usually pushes mask variables moving towards the positive infinity direction according to the significance of network connectivity. When training converges, both gradients on the mask variables will reach equilibrium at important network connections. Therefore, an optimized neural network with an expected level of sparsity is learned and determined through training.

F. Proposed SCL algorithm

Algorithm 1 describes the details of our proposed automatic sparse connectivity learning method. Through repeatedly calculating gradients of weight variables and mask variables, this algorithm updates them through stochastic gradient descent (SGD). When the network pruning process is complete, a sparse neural network is obtained. The well-trained sparse weights are calculated as

$$\mathbf{w}^* = \widetilde{\mathbf{w}} \odot H(\widetilde{\mathbf{m}}) \tag{22}$$

Algorithm 1 Automatic Sparse Connectivity Learning Input: A training dataset \mathcal{D} , a *L*-layer neural network, weight variables $\mathcal{W} = \{\widetilde{\mathbf{w}}^{(l)}\}_{l=1}^{l=L}$, mask variables $\mathcal{M} = \{\widetilde{\mathbf{m}}^{(l)}\}_{l=1}^{l=L}$, a connectivity decay coefficient λ_1 , an L_2 regularization coefficient λ_2 , and the total number of training epochs *T*.

Output: Well-trained sparse weights $W^* = \{\mathbf{w}^{(l)*}\}_{l=1}^{l=L}$. **Initialization:** Initial values of mask variables are positive. Initial values of weight variables are random according to [47].

1: repeat

2: Sample a mini-batch from \mathcal{D} as input data.

3: Forward pass:

First, calculate all the masks using Eq. (11). Next, compute all weights using Eq. (10). Then, Calculate all layers like a normal neural network with weights and sampled input.

4: Backward pass:

First, calculate weight gradients using Eq. (13). Next, relax mask gradients through the Identity STE using Eq. (19). Then, normalize mask gradients using Eq. (20).

5: Update weights:

Use gradients obtained from the backward pass to update all weight variables and mask variables via SGD.

- 6: **until** T epochs complete
- 7: Compute the well-trained sparse weights using Eq. (22) and output a sparse neural network.

G. Comparison of SCL with existing automatic pruning

It is necessary to comprehensively compare our SCL algorithm with existing automatic pruning methods in the literature (*i.e.*, [34]–[38]). Therefore, conceptual comparison and experimental results discussion of [34], [35], [37], [38] will be provided in Sections V-D and V-E. Since [36] does not report proper experimental results to compare, we only discuss their algorithm differences as below. Compared with [36], SCL has the following differences. First, the way of dealing with the non-differentiation problem of discrete masks is different. SCL uses a deterministic STE to estimate the gradient of discrete masks, whereas [36] uses the Gumbel-Softmax trick to relax the discrete masks to a continuous form. Second, the way to reduce the computational workload is different. [36] focuses on conditional computation using dynamic inference, so masks vary with network inputs. Besides, since the efficiency of data handling is greatly affected, the conditional computational graph in [36] is not friendly to hardware accelerators. In contrast, SCL focuses on network pruning that is static during inference, so the computational graph is fixed after training. Therefore, the feature of static computation reduction in SCL is more appropriate than [36] for hardware acceleration.

IV. BASELINE NETWORK ARCHITECTURES AND EXPERIMENTAL SETUPS

We conduct the experiments on four NVIDIA TITAN XP GPUs using PyTorch³. We choose VGGs [3], ResNets [4], DenseNets [5], and EfficientNets [48] as our baseline CNN architectures, because most of the state-of-the-art neural networks are based on them. ResNets and DenseNets are outstanding due to their high accuracy and fewer number of trainable parameters. EfficientNet is a lightweight high-accuracy convolutional neural network architecture with much fewer parameters. In CNN experiments, the proposed SCL technique is evaluated using the datasets of MNIST [49], CIFAR [50], and ImageNet (*i.e.*, full data of ImageNet2012 classification) [51]. In RNN experiments, SCL is evaluated on a language model using the WikiText-2 dataset [52].

For the experiments on the MNIST dataset, the baseline is a DenseNet-based network, where only fully connected layers are implemented (i.e., no convolutional networks involved) to evaluate the proposed SCL method. Then, two-dimension 28×28 image samples are flattened before being fed into the neural network. The input images are normalized using the channel mean and standard deviation. Within this baseline network, sixteen fully connected layers with a growth rate of 8 are applied to feature extraction, followed by a softmax function for object classification. For the experiments on the CIFAR dataset, convolutional networks are evaluated. The VGGs and ResNets are adapted from the codes⁴ of [53]. The codes of DenseNets are slightly different from the description in the paper of [5], which is referred to official codes⁵. For the experiments on the ImageNet dataset, we only evaluate SCL on VGG-16, ResNet-50⁶, and EfficientNet-B0⁷ due to the limitation of computational resources. The language model used in the RNN experiments consists of an encoder embedding, two LSTM layers, and a decoder embedding. As encoder and decoder embeddings are tied to improve the perplexity

⁷ https://github.com/lukemelas/EfficientNet-PyTorch

results [54], we only account for one of the encoder/decoder embeddings in sparsity statistics. The vocabulary size and LSTM hidden layer size are 33,278, and 1,500, respectively. The language model used in the RNN experiments is adapted from a PyTorch word language model⁸.

All the networks are trained by the SGD optimizer. We adopt the weight initialization method in [47], and batch normalization in [55] for fast training. The objective function of an L-layer network is expressed as

$$\mathcal{L} = \mathcal{C} + \lambda_1 \cdot \sum_{l=0}^{L} \sum_{i=0}^{|\mathbf{m}^{(l)}|} m_i^{(l)} + \lambda_2 \cdot \sum_{l=0}^{L} \sum_{i=0}^{|\widetilde{\mathbf{w}}^{(l)}|} (\widetilde{w}_i^{(l)})^2 \quad (23)$$

where C is the cross-entropy loss for classification, λ_1 and λ_2 are the coefficients of connectivity decay and L_2 regularization, respectively. $|\mathbf{m}^{(l)}|$ and $|\mathbf{\tilde{w}}^{(l)}|$ refer to the number of elements in $\mathbf{m}^{(l)}$ and $\mathbf{\tilde{w}}^{(l)}$, respectively. Note that λ_1 and λ_2 are not designer-defined pruning criteria or layer-wise pruning hyper-parameters (e.g., pruning threshold or ratio), which are indispensable in existing pruning methods in the literature. In our proposed SCL method, the layer-wise sparsity is automatically determined by the learning algorithm itself, without the intervention of designers. In our SCL method, tuning λ_1 can adjust the trade-off between network sparsity and accuracy. Therefore, designers can adjust λ_1 to reach the desired sparsity-accuracy balance. For a given value of λ_1 , we do not run full training iterations to check the resulting sparsity. Instead, we run a small portion (e.g., 10%) of the full training iterations to see if the target sparsity can be achieved. In this way, a proper λ_1 for the target sparsity can be found with several attempts. To further reduce the trial time, we start with a larger value of λ_1 , which leads to faster network pruning and hence shorter running time. Then, we decrease the value of λ_1 until the target sparsity is obtained. Therefore, the running time for tuning λ_1 is not huge. A similar process for tuning λ_1 is applied to other datasets and baseline architectures in the experiments of Section V. To make a concise comparison with existing works, we list the experimental results in terms of sparsity, the number of parameters, FLOPs, and accuracy in Tables II-XI. Note λ_2 is the coefficient of L_2 regularization, which has nothing with network pruning. L_2 regularization is required to ensure an effective learning rate when batch normalization is applied [56].

For experiments on the MNIST and CIFAR datasets, minibatch size and initial learning rate are set to 64, and 0.1, respectively. For experiments on the MNIST dataset, a simple training schedule is used, and mask variables do not update in the first and last 15 epochs. For experiments on the CIFAR datasets, mask variables do not update in the first 150 epochs to facilitate weight convergence. Then, mask variables are updated to obtain the corresponding sparse network connectivity. The learning rate is adjusted to 0.01 when the target sparsity is almost achieved. Finally, the mask variables continue to update for 20 epochs and achieve stable network connectivity. In order to achieve ultimate convergence stability, the weights continue to update for 80 epochs with a learning rate of

³ https://pytorch.org

⁴ https://github.com/Eric-mingjie/rethinking-network-pruning

⁵ https://github.com/liuzhuang13/DenseNet

⁶ https://github.com/pytorch/vision/tree/master/torchvision/models

⁸ https://github.com/pytorch/examples/tree/master/word_language_model

0.01 and then another 80 epochs with a learning rate of 0.001. Initial values of mask variables are positive. Weight variables are randomly initialized according to [47] and trained from scratch. A larger λ_1 usually leads to less number of training epochs for realizing the same sparsity expectation. For experiments on the ImageNet dataset, all the models are pre-trained before sparse training for saving time. Mini-batch size and initial learning rate are set to 128 and 0.005, respectively. We use 90 epochs for sparse connectivity training and another 60 epochs to achieve stable training procedure of Distiller⁹.

V. EXPERIMENTAL RESULTS AND COMPARISON

A. Effect of STE and mask gradient normalization

An map fitting experiment has been designed to evaluate the effect of STE and mask gradient normalization. A threelayer fully connected network is built as a baseline inputoutput mapping. Each network layer consists of 64 neurons, followed by batch normalization and ReLU activation. In order to ensure output convergence, weights are randomly initialized according to [47]. A set of inputs and mask variables are randomly selected and applied to this baseline. These inputs and mask variables follow a normal distribution with 0-mean and unit standard deviation. Thus, without any training, output results of the baseline are obtained as a benchmark.

Next, the baseline architecture is reused to perform output fitting experiments, where inputs and weights are the same as the baseline, but mask variables are randomly initialized and then trained through gradient descent to match the benchmark (*i.e.*, baseline outputs). These fitting experiments involve different STEs (*i.e.*, proxy mask gradients) for binary mask relaxation with or without mask gradient normalization. In these experiments, the ideal fitting result is that after training, mask variables finally converge and produce the same outputs as the benchmark, indicating an accurate reproduction of the baseline input-output mapping.

Figure 4 shows the mean squared errors (MSEs) obtained from the above experiments. MSE tells how close the outputs of trained networks are to the benchmark. The Clipped ReLU STE leads to the worst fitting results, because its gradient is 0 when a mask variable is negative or larger than a threshold as shown in Figure 2(b). The results of ReLU STE are better, because it has fewer dead zones as shown in Figure 2(a). Note that once a mask variable falls into a dead zone, there is no chance to get out. Compared with the Clipped ReLU and ReLU STEs, Leaky ReLU, Softplus, and Identity STEs achieve better and similar results. This observation is consistent with Eq. (16) and theoretical analysis in Section III-C. Positive proxy gradients of these three STEs guarantee a loss descent direction to push mask variables towards the minima. Mask gradient normalization reduces MSEs for all STEs, validating the necessity of normalizing mask gradients.





Fig. 4. Output fitting results with STEs for binary mask relaxation. Solid and dashed curves are with or without mask gradient normalization, respectively.

TABLE I TRAINED DENSENET-BASED NETWORKS ON MNIST.

Scheme	# Param.	Sparsity	Accuracy
Baseline	$117,\!152$	0 %	98.35%
$\overline{\lambda_1} = 0$	77,422	34.0%	98.47%
$\lambda_1 = 0.01$ $\lambda_1 = 0.03$	$10,153 \\ 4,488$	91.3% 96.2%	98.24% 98.01%
$\lambda_1 = 0.08$	1,375	98.8%	94.64%
$\lambda_1 = 0.1$	252	99.8%	78.46%

B. Sparse connectivity of fully connected network on MNIST

To quickly evaluate the effectiveness of the proposed SCL method, we carried out experiments using a DenseNet-based network with the MNIST database, as described in Section IV. As shown in Table I, a larger connectivity decay λ_1 corresponds to more sparse network connectivity. Even if λ_1 is set to 0, our experimental results show that the two gradients due to performance loss minimization and L_2 regularization (the first and third term in Eq. (23)) can push some mask variables to negative values, when their corresponding network connections do no contribute to accuracy. As a result, our SCL-induced neural network has a good sparsity of 34%, and its object classification accuracy of 98.47% outperforms the DenseNet-based baseline network (i.e., 98.35%). From a system accuracy perspective, this means that the best network connection for this example is inherently sparse. In addition, when λ_1 is 0.03, our SCL-induced sparse networks can achieve a very high accuracy of above 98% and sparsity of about 96.2%. Note that in this fully connected layer experiment, the sparsity is equal to the reduction in FLOPs. As a result, the SCL-induced DenseNet-based network achieves a 96.2% reduction in FLOPs (*i.e.*, approximately $26.3 \times$ lower than the baseline with an accuracy loss of only 0.34%). This experiment validates that our proposed SCL method supports effective network pruning on fully connected layer architectures.

Based on these experimental results, we add together the binary connections that are directly connected to the input data and then normalize them. The normalized connections for three typical λ_1 values are illustrated in Figure 5. It is clear that

⁹ https://github.com/NervanaSystems/distiller/tree/master/examples/word_l anguage_model



Fig. 5. Visualization of learned network connectivity on MNIST.



Fig. 6. Connection density profiles for all layers in DenseNet-40 and ResNet-110 for different overall densities.

network connections are mainly concentrated in central areas, so our SCL training ignores edges. As a result, the pixels of digits are trained to be connected in the central area. When λ_1 is set to 0.1 (*i.e.*, very high sparsity), only a few central pixels are connected as illustrated in Figure 5(a).

C. Connectivity learning analysis on CIFAR-10

In order to demonstrate the automatically learned network connectivity, we report and discuss the density profiles (i.e., the percentage of remaining network connections) of several SCLinduced networks on the CIFAR-10 dataset. Figure 6 shows that the proposed SCL method can automatically learn and determine the corresponding density profile for each network layer. In Figure 6(a), there are three layers with relatively high densities, which are the layers located before the three dense blocks. Since it is most often reused in this low-density DenseNet network, the three important network layers should retain more connections. In Figure 6(b), when this DenseNet network has a high targeted density of about 30%, most of the network layers obtained by our SCL method have a similar density. Note that the density of the last network layer is completely different in Figure 6(a) and Figure 6(c). This is due to different types of redundancy between DenseNet and ResNet structures. In DenseNet-40, the fact that previously extracted features are channel-wise concatenated to subsequent layers leads to 456 channel connections to the last layer. In ResNet-110, previously extracted features are element-wise added to subsequent layers, so the last layer contains only 64 channel connections and there is less redundancy. Therefore, DenseNet-40 removes most weights in the last layer, whose density is about 7% in Figure 6(a), while ResNet-110 retains most weights in the last layer, whose density is about 65% in Figure 6(c). In Figure 6(c) and 6(d), since ResNet-110 is too deep, network connections in some layers are completely removed. In Figure 6(d), even though 70% of weights and some computational expensive shallow layers are removed, our SCL-induced network is superior to the baseline networks (94.82% vs. 93.57% in Table III). In summary, we find that unlike conventional pruning methods which require designer-defined pruning criteria or hyper-parameter for each layer, our proposed SCL method can automatically learn and select important network connections for given baseline structures.

D. Comparison with the state-of-the-art pruning methods on CIFAR-10 and CIFAR-100

In order to comprehensively evaluate the proposed SCL method, we compare with the state-of-the-art pruning methods in the literature, including both non-structured and structured methods on the baseline networks of VGGs, ResNets, and DenseNets.

1) Comparison with non-structured pruning methods: We compare the proposed SCL with non-structured pruning methods, e.g., Frankle et al. [13] and Han et al. [7]. The results of Frankle et al. [13] (i.e., Lottery Ticket Hypothesis) are copied from Liu et al. [53]. As illustrated in Figure 7, we use a very high sparse target (97%) to learn sparse VGGs. Our SCL-induced VGGs is superior to the state-of-the-art pruning methods of [7], [13], [14], [20] in both sparsity and accuracy. For DenseNet-BC-100, with the same sparsity of 80%, our SCL method results in an accuracy of 95.40%, which is much higher than the accuracy of 95.04% in Han et al. [7]. Even though we don't compare FLOPs here due to we are unable to access pre-trained models in these previous works, the reduction in FLOPs is positively correlated with the sparsity for non-structured sparse models. Compared to the humandetermined criterion of non-structured pruning methods, the proposed task-aware SCL method yields better results.

2) Comparison with structured pruning methods: As shown in Table II and III, the ResNet networks trained by our proposed SCL method are more sparse and more accurate. The higher the sparsity, the more savings the FLOPs have. Besides, even with a higher sparsity (*e.g.*, 70% on ResNet-20 and 90% on ResNet-110), the SCL-induced ResNet neural networks achieve better accuracy than the baselines in [4]. In contrast, when the target sparsity exceeds about 30%, the existing structured pruning methods in [13] [14] [16] [57] exhibit significant accuracy degradation.

Table IV lists the experimental results for DenseNet neural networks on the CIFAR-10 dataset. When the target sparsity is moderate (*e.g.*, sparsity = 40% in Table IV), our SCL-induced networks show higher accuracy than the baseline in [5]. When setting a higher target sparsity (*e.g.*, 65%), the accuracy of Liu *et al.* [53] drops significantly. In contrast, even at a high sparsity of 90%, our SCL-induced DenseNet network



[7] and Liu et al. [20].

Fig. 7. Sparsity comparisons with non-structured pruning methods.

TABLE II RESULTS OF RESNET-20 ON CIFAR-10.

Scheme	# Param.	Sparsity	FLOPs \downarrow	Accuracy
Baseline [4]	$0.268 \mathrm{M}$	0%	0%	91.25%
He <i>et al.</i> [16] He <i>et al.</i> [16] He <i>et al.</i> [16]	0.241M 0.214M 0.188M	$10\% \\ 20\% \\ 30\%$	$15\% \\ 29\% \\ 42\%$	$92.24\% \\ 91.20\% \\ 90.83\%$
SCL SCL	0.133M 0.080M	$50\% \\ 70\%$	$49\% \\ 69\%$	92.61% 92.35%

TABLE III RESULTS OF RESNET-110 ON CIFAR-10.

Scheme	# Param.	Sparsity	FLOPs \downarrow	Accuracy
Baseline [4]	$1.72 \mathrm{M}$	0%	0%	93.57%
He et al. [16]	1.20M	30%	41%	93.86%
Yu et al. [57]	0.98M	43%	44%	93.39%
Yu et al. [57]	1.17 M	32%	39%	93.34%
Li et al. [14]	1.17M	32%	39%	93.36%
Frankle et al. [13]	$1.17 \mathrm{M}$	32%	39%	93.15%
SCL	$0.51 \mathrm{M}$	70%	73%	94.82%
SCL	$0.17 \mathrm{M}$	90%	90%	94.56%

TABLE IV **RESULTS OF DENSENET-40 ON CIFAR-10. RESULTS OF PREVIOUS** METHOD ARE COPIED FROM LIU et al. [53].

Scheme	# Param.	Sparsity	FLOPs \downarrow	Accuracy
Baseline [5]	$1.04\mathrm{M}$	0%	0%	94.76%
Liu et al. [20] Liu et al. [20]	$0.66\mathrm{M}$ $0.35\mathrm{M}$	$36\% \\ 65\%$	$28\% \\ 55\%$	94.81% 94.35%
SCL SCL SCL	0.62M 0.30M 0.10M	$40\% \\ 71\% \\ 90\%$	38% 70% 88%	94.81% 94.66% 94.53%

demonstrates little loss of accuracy. The superiority of our SCL is not only due to the advantage of the task-aware feature, but also because the resultant weight-level sparse models have a larger representative capacity.

3) Comparison with L_0 regularization method: L_0 regularization (Gumbel-Softmax trick) [34] determines zeroweight network connections by including a set of non-negative stochastic gates. Compared with [34], the SCL method has significant differences in the objective function, mask training method, actual reduction of sparse connectivity and FLOPs,



(c) DenseNet-BC-100, comparison with Han et al. [7].

as described below. First, the the objective function is different. In this work, the sparse regularization term in the objective function is the L_0 -norm of masks. In contrast, the regularization term in the objective function of [34] is a statistically expected L_0 -norm of masks. Second, the mask training method is different. Stochastic sampling is used in [34] to train mask variables, whereas this work directly trains mask variables without using stochastic sampling. Due to the use of stochastic sampling, the pruning method in [34] is complicated. In contrast, in this work, masks are constrained to 0 or 1 by applying a unit step function on it, and sparsity is produced by penalizing the L_0 -norm of weights. Hence, this work is simple and efficient. Third, the actual reduction of sparse connectivity and FLOPs is different. Both training and testing stages in this work are deterministic, whereas training and testing stages in [34] are stochastic and deterministic, respectively. Note a mask variable implies no connectivity only when its probability is zero. In this work, the probability of connectivity is either 0% or 100%, a reduction of L_0 -norm of masks (i.e., probability of certain connectivity is regulated from 100% to 0%) in the training stage indicates the same reduction in the testing stage. In contrast, the probability of connectivity is continuous from 0% to 100% in the training stage of [34]. Hence, in [34], the reduction of expected L_0 norm during training does not necessarily mean a reduction of L_0 -norm during testing. Besides, the discrete mask function is used in [34] for hard pruning in the testing stage. Discrete mask function potentially has a large discrepancy from the probabilistic continuous mask function during training time [35]. Furthermore, the use of stochastic sampling in [34] leads to a huge gap between the expected L_0 in the stochastic training phase and the actual L_0 in the deterministic testing phase. In [34], although the hard concrete distribution trick allows zero gates to be produced, the expected L_0 can not reflect the actual L_0 during inference.

Our proposed SCL takes advantage of STE [41] to redefine gradients of mask variables. Even though STE involves the use of STE in stochastic neurons, we think STE is also applicable to deterministic neurons, because deterministic neurons are regarded as a special case of stochastic neurons with a probability of either 0% (i.e., no connection) or 100% (i.e., with connection). The use of STE on deterministic problems has been empirically verified by the deep learning community [38], [45]. Then, L_0 norm of weights is integrated as a regularization

TABLE V Comparison with L_0 regularization [34] on CIFAR datasets. WRN-28-10 [58] is used as the baseline. "-" indicates results not reported.

Scheme	CIFA	AR-10	CIFAR-100		
Scheme	Sparsity Accurac		Sparsity Accuracy		
Baseline [58]	0%	96.00%	0%	80.75%	
Louizos <i>et al.</i> [34] Louizos <i>et al.</i> [34]	-	96.07% 96.17%	-	80.96% 81.25%	
SCL SCL SCL	$20\% \\ 51\% \\ 91\%$	96.36% 96.53% 96.33%	$20\% \\ 50\% \\ 90\%$	81.79% 81.87% 81.51%	

term of the objective function. Despite the role of SCL in training is similar to that of L_0 regularization, we think SCL is advantageous because its efficient processing of mask variables helps sparsity training. As a result, the proposed SCL method is more direct and effective in encouraging sparse network connections. Table V lists the experimental results of L_0 regularization [34] and SCL on the CIFAR-10 and CIFAR-100 datasets. The observation that both of them outperform the baseline model [58] indicates that the best network connections should be sparse. Moreover, the experimental results of SCL with three sparsity levels (i.e., 20%, 50%, and 90%) are provided. Their corresponding accuracy results are better than those of L_0 regularization, even though the sparsity results are not reported in [34]. In fact, when we use the hyper-parameters and codes in [34] to repeat the experiment for Table V, we find that the obtained network is not sparse. Although the expected L_0 norm of weights has been significantly reduced, it is still not low enough to generate sparse connections.

4) Comparison with existing STE-based pruning method: In addition to evaluating the effect of STE and mask gradient normalization in a mapping experiment in Section V-A, we also conduct experiments with DenseNet-based networks on MNIST and with ResNet-20 and ResNet-110 on CIFAR-10, respectively. We compare the results of SCL with the existing STE-based pruning method [38], which uses Leaky ReLU or Softplus STEs without mask gradient normalization. As shown in Tables VI and VII, under the same 95% sparsity, the result of SCL (*i.e.*, 90.08%) is superior to [38] (*i.e.*, 88.01% or 87.95%). Table VII also provides the comparison results of VGG-16 on CIFAR-10, as reported in [38].

TABLE VI Comparison Results of DenseNet-based networks on MNIST.

	Sparsity	Leaky ReLU	Softplus	Identity
w/o norm	95.0%	98.26% [38]	98.31% [38]	98.37%
w/ norm		98.33%	98.35%	98.39% (SCL)

E. Comparison with state-of-the-art pruning methods on ImageNet

In Tables VIII and IX, except for Han *et al.* [7], the performance results of existing pruning methods on ImageNet are copied from Lin *et al.* [18].

 TABLE VII

 COMPARISON OF RESNET-20, RESNET-110, AND VGG-16 ON CIFAR-10.

		Sparsity	Leaky ReLU	Softplus	Identity
ResNet-20	w/o norm w/ norm	95.0%	88.01% [38] 89.48%	87.95% [38] 90.13%	88.10% 90.08% (SCL)
ResNet-110	w/o norm w/ norm	95.0%	93.15% [38] 93.25%	92.26% [38] 93.44%	92.76% 93.58% (SCL)
VGG-16	w/o norm w/ norm	98.6% 97.0%	92.189	% [38]	94.39% (SCL)

In most convolutional neural networks, most of the weights are in fully connected layers, while most FLOPs are in convolutional layers. The number of fully connected layers for VGG-16 and ResNet is three and one, respectively. As shown in Table VIII, the structured pruning methods of [15], [18], [59] on VGG-16 lead to low sparsity and significant computational cost reduction, because they only prune the convolutional layers for fewer FLOPs and do not prune fully connected layers. On the other hand, the non-structured pruning methods, i.e., Han et al. [7] and our proposed SCL in this study, can obtain higher weight compression ratios, mainly due to the pruning of fully connected layers. Specifically, the method of [7] achieves a weight sparsity of 92.5% over its baseline. At a sparsity level of around 90%, even though the performance (i.e., TOP-1 and TOP-5 drop) of our proposed SCL is worse than [7], our SCL method can achieve a much higher accuracy of 69.20% over 68.66% in [7] and a more FLOPs reduction of 87% over 79% in [7]. The difference between our proposed SCL and [7] is as followed. The pruning efforts in [7] are mainly for fully connected layers, but less for pruning convolutional layers. In contrast, depending on the weight significance, the proposed SCL can efficiently perform pruning on both fully connected layers and convolutional layers. We also observe that the baseline of [7] may not converge, because its baseline TOP-1 accuracy is reported as 68.5%, and our TOP-1 accuracy is larger than 71.5%. The performance gain in [7] is attributed to the use of hundreds of training epochs during pruning, which leads to a much better convergence.

As shown in Table IX, for low-sparisty pruning, these cutting-edge pruning methods [14], [15], [18], [25], [37], [59] show significant performance degradation (*i.e.*, TOP-1, TOP-5) than our SCL-induced results. Furthermore, our SCL method achieves a high sparsity of 74% with a FLOPs reduction of 79%, while retaining a slight performance degradation.

Compared with [33], this work has three main differences. First, our SCL method is pruning criterion-free, whereas [33] intends to learn proper layer-wise pruning criterion from a set of designer-defined criteria. Due to the criterion-free nature, our SCL method does not need hyper-parameters in criterion. In contrast, [33] needs hyper-parameters to determine the number of filters to keep. Second, the pruning performance of [33] heavily depends on human experience, including designer-defined hyper-parameters and criteria set. In contrast, our SCL method automatically learns the optimized network connectivity in a task-aware manner. Third, binary masks and weight parameters are jointly updated in SCL, whereas mask and weight parameters are trained separately in [33]. As shown in Table IX, this work shows better network pruning results than [33]. SCL achieves a much lower drop of TOP-1 accuracy (*i.e.*, 0.23%) than [33] (*i.e.*, 1.69%) and a much lower drop of TOP-5 accuracy (*i.e.*, 0.40%) than [33] (*i.e.*, 0.83%), meanwhile largely reducing the number of FLOPs (*i.e.*, 79%) than [33] (*i.e.*, 61%).

Compared with [37], this work has three main differences. First, our SCL method uses binary masks to represent network connectivity, whereas [37] uses continuous scaling factors to represent network connectivity. In the proposed SCL method, no threshold is needed to train binary masks. In contrast, softthreshold needs to be trained to obtain non-negative scaling factors in [37]. In [37], if the value of a scaling factor is zero, it indicates no network connection. In [37], a positive scaling factor indicates the existence of network connection. Second, our SCL method uses scheduled SGD to train the binary masks, whereas [37] uses APG to train the scaling factor parameters. Third, [37] only prunes convolutional layers to reduce FLOPs, while our SCL method prunes convolutional layers and fully-connected layers to compress the weight size. As a result, the SCL method leads to much higher sparsity in network connectivity, and therefore fewer trainable parameters. As shown in Table IX, this work shows better network pruning performance than [37]. SCL achieves a much lower drop of TOP-1 accuracy (i.e., 0.23%) and a much lower drop of TOP-5 accuracy (*i.e.*, 0.40%) than [37], meanwhile reducing the number of parameters by 58% (i.e., 6.6M vs. 15.6M).

Even though this work and [38] are all inspired by the general concept of STE, the objective function optimization, update rule for mask parameters, and coarse gradient estimation used in this work are significantly different from [38]. In this work, the weight and mask parameters are updated in the same optimization iteration. Yet, the weight and mask parameters are updated separately in [38], which corresponds to high computational complexity. As a result, the training calculation cost of SCL is significantly reduced. In this work, the update rule is not based on assumptions, and the update rule for mask parameters is derived through the back-propagation algorithm. In contrast, gradients in [38] are modified by dividing true mask gradients (i.e., gradients obtained through back-propagation) by the absolute value of weights element-wisely. Two coarse gradient estimations (i.e., gradients of Leaky ReLU and Softplus function) are used in [38], whereas the straight-through gradient estimation (i.e., gradient of identity function) is used in this work. As shown in Table IX, in addition to the significant reduction of FLOPs (i.e., 79% in this work vs. 55% in [38]), our SCL pruning method achieves a much lower drop of TOP-1 accuracy (*i.e.*, 0.23% in this work vs. 0.40% in [38]). These experimental results demonstrate that our SCL method is better than [38].

Compared with [35], this work has two main differences. First, the SCP method in [35] assumes that feature maps follow a Gaussian distribution. Yet, this assumption is too strict to derive accurate gradients. In contrast, our SCL method does not rely on any assumptions. Our SCL method trains network connectivity through STE gradient estimation and gradient back-propagation. Second, the SCP method in [35] has a good pruning performance in network layers that are followed by BN and ReLU. However, there are no BN and ReLU in many deep neural networks. For example, BN does not exist in the VGGs network. In the architectures of MobileNets or EfficientNets, some convolutional layers are only followed by BN rather than BN and ReLU. As shown in Table 5 of [35], the network pruning performance is significantly degraded if only BN exists for channel pruning. As a result, the SCP method in [35] does not show good pruning performance for many neural networks. In contrast, experimental results in this work show that our SCL method is applicable to various neural networks, including VGGs, DenseNets, ResNets, EfficientNets, and RNNs. As shown in Table IX, SCL shows better network pruning performance than [35]. SCL achieves a much lower drop of TOP-1 accuracy (*i.e.*, 0.23%) than [35] (*i.e.*, 1.69%) and a much lower drop of TOP-5 accuracy (i.e., 0.40%) than [35] (*i.e.*, 0.98%), meanwhile largely reducing the number of FLOPs (i.e., 79%) than [35] (i.e., 54%).

We also apply SCL to the EfficientNet [48] architecture to learn sparse connections. Unlike VGG, ResNet, and DenseNet architectures that are designed by humans, the EfficientNet architecture is automatically determined by a neural architecture search technique. Depth-wise convolution is widely used in EfficientNet architectures to improve efficiency. As shown in Table X, three pruning polices are used by us for depth-wise convolution in the experiments of [11] on the ImageNet dataset. In the first pruning policy, the pruning rates of the convolution, depth-wise convolution, and classifier layers are set to 0.2, 0.1, and 0.2, respectively. In the second pruning policy, the pruning rates of the convolution, depthwise convolution, and classifier layers are set to 0.5, 0.5, and 0.5, respectively. In the third pruning policy, the pruning rates of the convolution, depth-wise convolution, and classifier layers are set to 0.5, 0.125, and 0.6, respectively. Compared with the baseline EfficientNet-B0, Table X demonstrates that when the sparsity is low (such as 19.5% and 25.8%), both the pruning method of Zhu *et al.* [11] and SCL show a negligible accuracy degradation. According to the magnitude of weights, the work of Zhu et al. [11] forces smaller weights to zero. As a result, although it is an element-wise pruning method, if the difference of weight magnitudes between depth-wise channels is large, the pruning method of Zhu et al. [11] tends to completely remove certain depth-wise convolution channels. The second policy of Zhu et al. [11] sets the weights of all layers to be pruned by 50%, experimental results show that FLOPs are reduced by 60% with a big TOP-1 accuracy drop of 0.79%. The third policy of Zhu et al. [11] prunes less weights of depth-wise layers and obtains a better accuracy. These pruning results of Zhu *et al.* [11] show that it is difficult to find an appropriate pruning policy to obtain satisfactory pruning results. In contrast, since the network connectivity learned by SCL is determined by the significance of weight, SCL automatically retains necessary depth-wise channels even if their magnitudes of weight are small. As a result, when SCL prunes more weights than the pruning method of Zhu et al. [11] (*i.e.*, a sparsity of 55.0% versus 50.9%), SCL significantly improves the accuracy results. Thus, SCL also outperforms the pruning method of Zhu et al. [11] in terms of EfficientNet architectures.

F. Sparse RNN learning on WikiText-2

When SCL prunes weight matrices in RNNs (i.e., setting some weight values to zero), it does not necessarily produce zero gradients. In addition, the network connectivity of RNNs is not sparsely initialized in SCL. During the training process, sparse connections are gradually produced by SCL. Hence, SCL does not cause the exploding/vanishing gradient problem for RNNs. In this work, SCL is applied to a word-level language model WikiText-2 [52] for verifying its effectiveness on RNNs. The perplexity of the language model is evaluated. The lower the perplexity, the better the RNN model. Table XI lists the experimental results of existing pruning methods (i.e., Zhu [11] and Narang [60]) and SCL. To obtain these results, the pruning method of Zhu et al. [11] has to find an appropriate combination of sparse rates for each weight tensor through a lot of trials and errors, while the sparsity for each weight is automatically found by SCL. Compared with the baseline model that has a sparsity of 0%, sparse models obtained by these existing state-of-the-art pruning methods and SCL achieve higher perplexity values. The 80% sparse RNN model trained by Zhu et al. [11] or SCL reduces the number of weight parameters by almost 80% and meanwhile outperforms the baseline in terms of perplexity on the test dataset. The results indicate that SCL can effectively output excellent sparse connectivity for RNNs. Besides, when the expected sparsity is 95%, the perplexity results of SCL are better than those of Zhu et al. [11] and Narang [60].

G. Applicability of SCL in Sparse IndRNN learning

IndRNN [61] has recently been proposed to solve the exploding/vanishing gradient problem in RNNs. IndRNN modifies and updates the RNN states from

$$\mathbf{h}_t = \sigma(\mathbf{W}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{b}) \tag{24}$$

$$\mathbf{h}_t = \sigma(\mathbf{W}\mathbf{x}_t + \mathbf{u} \odot \mathbf{h}_{t-1} + \mathbf{b}) \tag{25}$$

where $\mathbf{x}_t \in \mathbb{R}^M$ and $\mathbf{h}_t \in \mathbb{R}^N$ are the input states and hidden states at a time step t, respectively. $\mathbf{W} \in \mathbb{R}^{N \times M}$, $\mathbf{U} \in \mathbb{R}^{N \times N}$, and $\mathbf{b} \in \mathbb{R}^N$ represent the weights for the current input, recurrent input, and bias of neurons, respectively. $\mathbf{u} \in \mathbb{R}^N$ is a weight vector. When the weight matrix \mathbf{U} happens to be a diagonal matrix, the vector \mathbf{u} can be regarded as the diagonal vector of matrix \mathbf{U} . N represents the number of neurons in this layer and σ represents an activation function. We can see that $\mathbf{U}\mathbf{h}_{t-1}$ in the RNN is reformulated as Hadamard product $\mathbf{u} \odot \mathbf{h}_{t-1}$ in the IndRNN. Therefore, the gradient of the *n*-th neuron at the time step t is changed from the RNN gradient

$$\frac{\partial J}{\partial h_t} = \frac{\partial J}{\partial h_T} \prod_{k=t}^{T-1} diag(\sigma'(h_{k+1})) \mathbf{U}^T$$
(26)

to the IndRNN gradient

$$\frac{\partial J_n}{\partial h_{n,t}} = \frac{\partial J_n}{\partial h_{n,T}} u_n^{T-t} \prod_{k=t}^{I-1} \sigma'_{n,k+1}$$
(27)

where $diag(\sigma'(h_{k+1}))$ is the Jacobian matrix of the element-wise activation function. Due to the term of $\prod_{k=t}^{T-1} diag(\sigma'(h_{k+1}))\mathbf{U}^T$, it is difficult to control the gradient of a RNN within a appropriate range. Fortunately, the exploding/vanishing gradient problem is easily addressed in the IndRNN by regulating the exponential term of $u_n^{T-t}\prod_{k=t}^{T-1}\sigma'_{n,k+1}$ within an appropriate range during the training process.

Next, let us analyze the impact of our SCL method on the exploding/vanishing gradient problem in IndRNN neural networks. For IndRNN networks, the proposed SCL method is applicable to prune weights in the matrix W in the IndRNN without affecting the exploding/vanishing gradient problem in training. Moreover, the weights in the vector **u** should be excluded from being pruned by SCL, because if some weights in **u** are pruned, some values of u_n^{T-t} will be set to zero, thus causing the vanishing gradient problem in the IndRNN. In fact, by comparing the dimension of weight matrix $\mathbf{W} \in \mathbb{R}^{N \times M}$ and weight vector $\mathbf{u} \in \mathbb{R}^N$, we see that the majority of weight parameters are located in the matrix W. As a result, when SCL prunes the weight matrix U and meanwhile prevents the vector **u** from being pruned, the pruning ability and space in the IndRNN networks are not significantly reduced. Since we have shown the effectiveness of the proposed SCL on RNNs, both of whose weight matrices W and U are pruned by SCL, it is expected that when only pruning the weight matrix W, the proposed SCL method can achieve highly sparse IndRNN networks without affecting the exploding/vanishing gradient problem.

VI. CONCLUSION

We present a Sparse Connectivity Learning (SCL) method to automatically explore and optimize sparse network connectivity. As the number of neural network connections is incorporated into the objective function, the network connectivity can be optimized for a given sparsity expectation to achieve the best performance. Our proposed SCL method has the task-aware ability, which does not require designer-defined pruning criteria or hyper-parameters for each network layer. As a result, the SCL-induced sparse networks are explored in a larger hypothesis space, and they have the potential to generate optimized network connections. The proposed SCL is applicable to various neural network architectures including fully connected networks, convolutional neural networks (VGGs, ResNets, DenseNets, and EfficientNets), and recurrent neural networks (RNNs). Experiments on the MNIST, CIFAR-10, CIFAR-100, ImageNet, and WikiText-2 datasets demonstrate that the proposed SCL method achieves highly efficient learning of sparse network connectivity in network compression ratio, FLOP reduction, and accuracy over existing state-of-the-art pruning methods in the literature. So far, SCL supports convolutions, fully connected layers, and RNN layers. We will explore SCL on other operators in the future.

ACKNOWLEDGEMENT

This work is partially supported by China Natural Science Foundation under grant (No. 62171391), the Opening Foun-

Scheme	# Param.	Sparsity	FLOPs \downarrow	$\text{TOP-1} \downarrow (\text{TOP-1})$	$\text{TOP-5} \downarrow (\text{TOP-5})$
Li <i>et al.</i> [14]	126.7M	8.38%	71%	0.29%	-0.05%
Luo <i>et al.</i> [15]	131.5M	4.92%	68%	-1.46%	-1.09%
Hu <i>et al.</i> [59]	126.7M	8.38%	71%	-0.64%	-0.43%
Lin <i>et al.</i> [18]	126.2M	8.75%	71%	-1.65%	-0.97%
Han <i>et al.</i> [7]	10.3M	92.5%	79%	-0.26% (68.66%)	0.44% (89.12%)
SCL	60.2M	56.5%	50%	-1.11% (72.84%)	-0.54% (90.88%)
SCL	36.9M	73.3%	71%	-0.33% (72.05%)	-0.25% (90.60%)
SCL	13.9M	89.9%	87%	2.49% (69.24%)	1.11% (89.24%)

TABLE VIII Pruning results of VGG-16 on ImageNet.

TABLE IXPruning results of ResNet-50 on ImageNet.

Scheme	# Param.	Sparsity	$FLOPs\downarrow$	TOP-1 ↓	TOP-5 \downarrow
Li et al. [14]	$15.9 \mathrm{M}$	38%	54%	3.36%	2.08%
Li et al. [14]	12.2M	52%	59%	4.31%	2.42%
Luo et al. [15]	$16.9 \mathrm{M}$	34%	41%	3.09%	1.63%
Luo et al. [15]	12.3M	52%	59%	4.12%	2.28%
Wen et al. [25]	13.2M	48%	49%	4.58%	2.68%
Hu et al. [59]	$15.9 \mathrm{M}$	38%	54%	3.47%	2.39%
Hu et al. [59]	12.2M	52%	59%	4.25%	2.41%
Lin et al. [18]	15.5M	39%	54%	2.83%	1.57%
Lin et al. [18]	12.0M	53%	59%	3.65%	2.11%
He et al. [33]	-	-	61%	1.69%	0.83%
Huang et al. [37]	15.6M	39%	43%	4.30%	2.07%
Xiao et al. [38]	-	-	55%	0.40%	-
Kang <i>et al.</i> [35]	-	-	54%	1.69%	0.98%
SCL	$17.9 \mathrm{M}$	30%	24%	-0.30%	-0.16%
SCL	6.6M	74%	79%	0.23%	0.40%

 TABLE X

 PRUNING RESULTS OF EFFICIENTNET-B0 [48] ON IMAGENET. ¹, ², AND ³

 INDICATE THE FIRST, SECOND, AND THIRD PRUNING POLICY,

 RESPECTIVELY.

Scheme	# Param.	Sparsity	FLOPs \downarrow	TOP-1 ↓	TOP-5 \downarrow
Zhu <i>et al.</i> [11] ¹ Zhu <i>et al.</i> [11] ² Zhu <i>et al.</i> [11] ³	4.24M 2.65M 2.59M	$19.5\%\ 49.7\%\ 50.9\%$	$18\% \\ 60\% \\ 51\%$	$\begin{array}{c} 0.04\% \\ 0.79\% \\ 0.41\% \end{array}$	$-0.04\% \\ 0.26\% \\ 0.18\%$
SCL SCL	$3.91\mathrm{M}$ $2.37\mathrm{M}$	$25.8\% \\ 55.0\%$	$24\% \\ 53\%$	$\begin{array}{c} 0.02\% \\ 0.32\% \end{array}$	$\begin{array}{c} 0.01\% \\ 0.13\% \end{array}$

 TABLE XI

 Results of sparse RNN learning on WikiText-2 [52]. Perplexity

 of the language models is evaluated, the lower perplexity the better.

Scheme	# Param.	Sparsity	Perpl. (Validation)	Perpl. (Test)
Baseline	$85.9 \mathrm{M}$	0 %	87.49	83.85
Zhu et al. [11] Zhu et al. [11] Narang et al. [60] Zhu et al. [11]	16.8M 8.6M 4.9M 4.3M	80.4% 90.0% 94.3% 95.0%	$89.31 \\90.70 \\100.23 \\98.42$	83.64 85.67 95.35 92.79
SCL SCL	$\begin{array}{c} 16.1\mathrm{M} \\ 4.6\mathrm{M} \end{array}$	81.3% 94.7%	88.97 97.63	$83.16 \\ 91.27$

dation of Yulin Research Institute of Big Data (Grant No. 2020YJKY04).

We thank Dr. Chun-Chen Liu from Kneron Inc. who provided insight and expertise that greatly assisted the research. We would also like to acknowledge Cheng-Hung Hsieh, Jia-En Hsieh from National Chiao Tung University for their helpful discussion.

REFERENCES

- [1] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.
- [2] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [3] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *ICLR*, 2015.
- [4] 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, 2016, pp. 770–778.
- [5] 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*, 2017, pp. 4700–4708.
- [6] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, p. 436, 2015.
- [7] S. Han, J. Pool, J. Tran, and W. Dally, "Learning both weights and connections for efficient neural network," in Advances in neural information processing systems, 2015, pp. 1135–1143.
- [8] Y. Guo, A. Yao, and Y. Chen, "Dynamic network surgery for efficient dnns," in Advances In Neural Information Processing Systems, 2016, pp. 1379–1387.
- [9] K. Ullrich, E. Meeds, and M. Welling, "Soft weight-sharing for neural network compression," in *ICLR*, 2017.
- [10] D. Molchanov, A. Ashukha, and D. Vetrov, "Variational dropout sparsifies deep neural networks," in *ICML*, 2017.
- [11] M. Zhu and S. Gupta, "To prune, or not to prune: exploring the efficacy of pruning for model compression," *arXiv preprint arXiv:1710.01878*, 2017.
- [12] E. Tartaglione, S. Lepsøy, A. Fiandrotti, and G. Francini, "Learning sparse neural networks via sensitivity-driven regularization," in *Advances* in *Neural Information Processing Systems*, 2018, pp. 3878–3888.
- [13] J. Frankle and M. Carbin, "The lottery ticket hypothesis: Finding sparse, trainable neural networks," in *ICLR*, 2019.
- [14] H. Li, A. Kadav, I. Durdanovic, H. Samet, and H. P. Graf, "Pruning filters for efficient convnets," in *ICLR*, 2017.
- [15] J.-H. Luo, J. Wu, and W. Lin, "Thinet: A filter level pruning method for deep neural network compression," in *ICCV*, 2017, pp. 5058–5066.
- [16] Y. He, G. Kang, X. Dong, Y. Fu, and Y. Yang, "Soft filter pruning for accelerating deep convolutional neural networks," *IJCAI*, 2018.
- [17] S. Lin, R. Ji, Y. Li, Y. Wu, F. Huang, and B. Zhang, "Accelerating convolutional networks via global & dynamic filter pruning." in *IJCAI*, 2018, pp. 2425–2432.
- [18] S. Lin, R. Ji, Y. Li, C. Deng, and X. Li, "Toward compact convnets via structure-sparsity regularized filter pruning," *IEEE transactions on neural networks and learning systems*, 2019.
- [19] Y. He, P. Liu, Z. Wang, Z. Hu, and Y. Yang, "Filter pruning via geometric median for deep convolutional neural networks acceleration," in *CVPR*, 2019.
- [20] Z. Liu, J. Li, Z. Shen, G. Huang, S. Yan, and C. Zhang, "Learning efficient convolutional networks through network slimming," in *ICCV*, 2017, pp. 2755–2763.

- [21] Y. He, X. Zhang, and J. Sun, "Channel pruning for accelerating very deep neural networks," in *ICCV*, 2017.
- [22] Z. Zhuang, M. Tan, B. Zhuang, J. Liu, Y. Guo, Q. Wu, J. Huang, and J. Zhu, "Discrimination-aware channel pruning for deep neural networks," in *Advances in Neural Information Processing Systems*, 2018, pp. 875–886.
- [23] C. Liu and H. Wu, "Channel pruning based on mean gradient for accelerating convolutional neural networks," *Signal Processing*, vol. 156, pp. 84–91, 2019.
- [24] J. Li, Q. Qi, J. Wang, C. Ge, Y. Li, Z. Yue, and H. Sun, "Oicsr: Outin-channel sparsity regularization for compact deep neural networks," in *CVPR*, 2019.
- [25] W. Wen, C. Wu, Y. Wang, Y. Chen, and H. Li, "Learning structured sparsity in deep neural networks," in *Advances in Neural Information Processing Systems*, 2016, pp. 2074–2082.
- [26] S. Anwar, K. Hwang, and W. Sung, "Structured pruning of deep convolutional neural networks," ACM Journal on Emerging Technologies in Computing Systems, vol. 13, no. 3, p. 32, 2017.
- [27] 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, 2018, pp. 6848–6856.
- [28] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
- [29] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1492– 1500.
- [30] J. Ye, X. Lu, Z. Lin, and J. Z. Wang, "Rethinking the smaller-normless-informative assumption in channel pruning of convolution layers," *arXiv preprint arXiv*:1802.00124, 2018.
- [31] P. Molchanov, S. Tyree, T. Karras, T. Aila, and J. Kautz, "Pruning convolutional neural networks for resource efficient inference," in *ICLR*, 2017.
- [32] L. Theis, I. Korshunova, A. Tejani, and F. Huszar, "Faster gaze prediction with dense networks and fisher pruning," arXiv preprint arXiv:1801.05787, 2018.
- [33] Y. He, Y. Ding, P. Liu, L. Zhu, H. Zhang, and Y. Yang, "Learning filter pruning criteria for deep convolutional neural networks acceleration," in *CVPR*, 2020, pp. 2009–2018.
- [34] C. Louizos, M. Welling, and D. P. Kingma, "Learning sparse neural networks through l_0 regularization," in *ICLR*, 2018.
- [35] M. Kang and B. Han, "Operation-aware soft channel pruning using differentiable masks," in *International Conference on Machine Learning*. PMLR, 2020, pp. 5122–5131.
- [36] C. Herrmann, R. S. Bowen, and R. Zabih, "Channel selection using gumbel softmax," in *European Conference on Computer Vision*. Springer, 2020, pp. 241–257.
- [37] Z. Huang and N. Wang, "Data-driven sparse structure selection for deep neural networks," in ECCV, 2018, pp. 304–320.
- [38] X. Xiao, Z. Wang, and S. Rajasekaran, "Autoprune: Automatic network pruning by regularizing auxiliary parameters," *Advances in neural information processing systems*, vol. 32, 2019.
- [39] E. Jang, S. Gu, and B. Poole, "Categorical reparameterization with gumbel-softmax," in *ICLR*, 2017.
- [40] C. J. Maddison, A. Mnih, and Y. W. Teh, "The concrete distribution: A continuous relaxation of discrete random variables," in *ICLR*, 2017.
- [41] Y. Bengio, N. Léonard, and A. Courville, "Estimating or propagating gradients through stochastic neurons for conditional computation," arXiv preprint arXiv:1308.3432, 2013.
- [42] M. Courbariaux, I. Hubara, D. Soudry, R. El-Yaniv, and Y. Bengio, "Binarized neural networks: Training deep neural networks with weights and activations constrained to+ 1 or-1," *arXiv preprint arXiv:1602.02830*, 2016.
- [43] P. Yin, S. Zhang, J. Lyu, S. Osher, Y. Qi, and J. Xin, "Blended coarse gradient descent for full quantization of deep neural networks," *Research in the Mathematical Sciences*, vol. 6, no. 1, pp. 1–23, 2019.
- [44] Z. Cai, X. He, J. Sun, and N. Vasconcelos, "Deep learning with low precision by half-wave gaussian quantization," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5918–5926.
- [45] P. Yin, J. Lyu, S. Zhang, S. Osher, Y. Qi, and J. Xin, "Understanding straight-through estimator in training activation quantized neural nets," *arXiv preprint arXiv:1903.05662*, 2019.

- [46] G. Hinton, "Neural networks for machine learning coursera video lectures," 2012.
- [47] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1026–1034.
- [48] M. Tan and Q. V. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *ICLR*, 2019.
- [49] Y. LeCun, C. Corinna, and B. Christopher J.C., "The mnist database of handwritten digits," http://yann.lecun.com/exdb/mnist/index.html.
- [50] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," University of Toronto, Tech. Rep., 2009.
- [51] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge," *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [52] S. Merity, C. Xiong, J. Bradbury, and R. Socher, "Pointer sentinel mixture models," in *ICLR*, 2017.
- [53] Z. Liu, M. Sun, T. Zhou, G. Huang, and T. Darrell, "Rethinking the value of network pruning," *ICLR*, 2019.
- [54] H. Inan, K. Khosravi, and R. Socher, "Tying word vectors and word classifiers: A loss framework for language modeling," in *ICLR*, 2017.
- [55] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.
- [56] T. van Laarhoven, "L2 regularization versus batch and weight normalization," in Advances in neural information processing systems, 2017.
- [57] R. Yu, A. Li, C.-F. Chen, J.-H. Lai, V. I. Morariu, X. Han, M. Gao, C.-Y. Lin, and L. S. Davis, "Nisp: Pruning networks using neuron importance score propagation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 9194–9203.
- [58] S. Zagoruyko and N. Komodakis, "Wide residual networks," in *BMVC*, 2016.
- [59] H. Hu, R. Peng, Y.-W. Tai, and C.-K. Tang, "Network trimming: A datadriven neuron pruning approach towards efficient deep architectures," *arXiv preprint arXiv*:1607.03250, 2016.
- [60] S. Narang, G. Diamos, S. Sengupta, and E. Elsen, "Exploring sparsity in recurrent neural networks," in *ICLR*, 2017.
- [61] S. Li, W. Li, C. Cook, C. Zhu, and Y. Gao, "Independently recurrent neural network (indrnn): Building a longer and deeper rnn," in *CVPR*, 2018, pp. 5457–5466.



Zhimin Tang received B.S. degree from Hunan Normal University in 2012 and M.S. degree from Hunan University in 2015. He is currently pursuing the Ph.D. degree in the Department of Automation at Xiamen University, China. He was a visiting scholar at the Department of Electrical and Computer Engineering of Southern Illinois University Carbondale from 2018 to 2020. His research interests include image processing, machine learning, and computer vision.



Linkai Luo is a Professor in the Department of Automation, Xiamen University, China. His research interests include machine learning, pattern recognition, data mining, computer vision, biological information processing and financial data analysis.



Bike Xie is currently a director of engineering in Kneron Inc, San Diego. He is and has been working on system architecture design for AI chip, deep learning model compression, and algorithm design for deep learning applications. He received his B.S. degree in electronic engineering from Tsinghua University, Beijing, in 2005. He received his M.S. and the Ph.D. degrees in electrical engineering from the University of California, Los Angeles in 2006 and 2010 respectively. At UCLA, he worked on a broad range of research topics including capacity regions

and encoding schemes for broadcast channels, download-time regions for peer-to-peer networks, packet coding and exchange for broadcast networks, universal turbo codes for space-time channels, and channel code design for optical communications. Dr. Xie then joined Marvell Semiconductor Inc., Santa Clara in 2010. At Marvell, he led a system team design physical layer systems and specifications for 6Gbps to 112Gbps SerDes systems, and capacitive touchscreen systems. In 2017, he joined Kneron Inc.



Yiyu Zhu received B.S. degree and M.S degree in Electrical Engineering from University of California, San Diego in 2017 and 2018. He worked at Kneron Inc as an algorithm engineer between 2018 and 2021. During this time, he developed algorithm for model compression and inference acceleration. His research interests include machine learning, computer architecture, and high performance computing.



Rujie Zhao received a B.S. degree in Electrical Engineering from Jiang Su Normal University and an M.S. degree from the State University of New York at New Paltz. Since fall 2018, he is working on his Ph.D. degree in the Department of Electrical and Computer Engineering at Southern Illinois University Carbondale, Carbondale, IL, USA. His research interests include deep neural network algorithms and architectures, memristor-based neural network implementation and optimization.



Lvqing Bi received the B.S., M.S. degrees in Electronic Science and technology from the Guangxi University, Nanning, China, in 2003 and 2013, respectively. Since 2015, he has been working toward the Ph.D. degree in Electronic Science and technology at the School of Electronic Science and Engineering, Xiamen University, Xiamen, China. His current research interests include artificial neural networks, intelligent manufacturing technology, network signal processing, terahertz metamaterials, information entropy and fuzzy sets.



Chao Lu received B.S. degree in electrical engineering from Nankai University and M.S. degree from the Hong Kong University of Science and Technology. He obtained Ph.D. degree at Purdue University, West Lafayette in 2012. From 2013 to 2015, he worked in US industry. Since July 2020, he became an associate professor of the Electrical and Computer Engineering Department of Southern Illinois University, Carbondale, IL, USA. His research interests include the design of high-performance circuit and system design, and deep machine learning

algorithms. Mr. Lu was the recipient of the Best Paper Award of the International Symposium on Low Power Electronics and Design in 2007, Best Paper Award Nomination of IEEE System-on-Chip Conference (2016), and Top 5 team in the International Hardware Design Contest of IEEE Design Automation Conference (2017).