

Restricted Random Pruning at Initialization for High Compression Range

Anonymous authors

Paper under double-blind review

Abstract

Pruning at Initialization (PaI) makes training overparameterized neural networks more efficient by reducing the overall computational cost from training to inference. Recent PaI studies showed that random pruning is more effective than ranking-based pruning, which learns connectivity. However, the effectiveness of each pruning method depends on the existence of skip connections and the compression ratio (the before-after pruning parameter ratio). While random pruning performs better than ranking-based pruning on architectures with skip connections, the superiority without skip connections is reversed in the high compression range. This paper proposes Minimum Connection Assurance (MiCA) that achieves higher accuracy than conventional PaI methods for architectures with and without skip connections, regardless of the compression ratio. MiCA preserves the random connection between the layers and maintains the performance at high compression ratios without the costly connection learning that ranking-based pruning requires. Experiments on CIFAR-10 and CIFAR-100 show that MiCA enhances the compression ratio and accuracy trade-offs compared to existing PaI methods. In VGG-16 with CIFAR-10, MiCA improves the accuracy of random pruning by 27.0% at $10^{4.7} \times$ compression ratio. Furthermore, experimental analysis reveals that increasing the utilization of the nodes through which information flows from the first layer is essential for maintaining high performance at a high compression ratio.

1 Introduction

Although deep neural networks (DNNs) have high generalization capability, both their training and inference are computationally expensive (Arora et al., 2019; Zhang et al., 2019; 2021; Neyshabur et al., 2019; Wen et al., 2022). These high costs arise because their computation depends on a large amount of parameters (Shoeybi et al., 2019; Brown et al., 2020; Dosovitskiy et al., 2021; Woo et al., 2023).

Network pruning achieves high generalization capability despite fewer parameters and can solve this problem. There are various types of pruning, including methods that train while sparsifying the network gradually by penalty terms (Chauvin, 1988; Weigend et al., 1990; Ishikawa, 1996), prune the network after training and then finetune it (LeCun et al., 1989; Hassibi et al., 1993; Lee et al., 2021), and prune and learn iteratively (Frankle & Carbin, 2019; Frankle et al., 2019; Renda et al., 2020). However, these aim to reduce the inference computational cost and need to train the dense model. By contrast, dynamic sparse training (Mocanu et al., 2018; Evci et al., 2020; Jayakumar et al., 2020) and pruning at initialization (PaI) train with sparse networks, thus reducing training costs. In particular, PaI has the lowest training computational cost among pruning methods because the network structure is fixed (Price & Tanner, 2021).

Basically, PaI calculates a criterion to determine which parameters are essential and selects the parameters to be pruned based on it. This type of PaI called *ranking-based pruning at initialization* (RbPI) (Lee et al., 2019; Wang et al., 2020; Tanaka et al., 2020) can learn the network connections explicitly but needs to calculate the criterion using an expensive process such as backpropagation. On the other hand, another type of PaI called *random pruning at initialization* (RPI) has a negligibly small additional cost because it only prunes a network randomly without calculating a criterion. At first glance, RbPI seems to perform better than RPI since it learns connections, but some works suggested that RPI could construct subnetworks

with similar or better performance obtained by RbPI. Frankle et al. (2021) revealed that RPI and RbPI had comparable accuracy at $1-10^2 \times$ compression ratios when applying the same sparsity set separately for each layer (i.e., the sparsity distribution). Similarly, the work by Su et al. (2020) showed that RPI with ad-hoc sparsity distribution improved the trade-off between parameter ratio of dense to sparse network—*compression ratio*—and accuracy than RbPI. Furthermore, randomly pruned networks outperform dense networks in aspects such as out-of-distribution detection and adversarial robustness (Liu et al., 2022). Thus, RPI seems to combine simple pruning processing with high performance among PaI methods.

On the other hand, a recent thread of PaI research (Vysogorets & Kempe, 2023) showed a curious phenomenon of RPI: its efficiency at high compression ratios depends on skip connections in the DNN. At more than $10^2 \times$ compression ratios, randomly pruned models without skip connections are less accurate than RbPI. It differs from the result in low compression ratio by Frankle et al. (2021) and implicates that highly sparse networks need to learn connection. However, whether connection learning is essential in the high compression range is debatable. For instance, the work by Gadhikar et al. (2023) improved the performance of skip connection-free architecture by adding parameters (i.e., edges) to non-functional neurons in a randomly pruned network at $10-10^3 \times$ compression ratios. It indicates that even random pruning can improve performance if the connections between layers are preserved. However, this approach is not essential for higher compression ratios because the additional edges inhibit compression. Therefore, there is a need for more essential solutions to improve performance in skip connection-free architecture.

In order to address this problem, this paper introduces a novel PaI algorithm for high compression range: *Minimum Connection Assurance* (MiCA). Specifically, it preserves top-to-bottom information propagation among randomly pruned layers by building a random connection—*minimum connection*—using some of the pre-allocated edges. The minimum connection is constructed by pre-determining and connecting the neurons that the subnetwork uses, and the subnetwork with the connection maintains the pre-defined sparsity distribution even when connecting its neurons randomly. Thus, all allocated edges can be functional even in a high compression range. Since MiCA has this small constraint on the placement of the edges while keeping the connection random, it stands as *restricted* RPI algorithm in the field of PaI algorithms. We evaluate MiCA on CIFAR-10 and CIFAR-100 datasets (Krizhevsky et al., 2009) with VGG (Simonyan & Zisserman, 2014) and ResNet (He et al., 2016) architectures. MiCA enhances the performance for not only skip connection-free architectures such as VGG but also architectures with skip connections such as ResNet. Furthermore, despite the random connection, MiCA improves the trade-off between compression ratio and inference accuracy compared to RPI and RbPI methods. In other words, MiCA shows that the connections learned by RbPI can be replaced by random connections even in the high-compression range.

The rest of the paper is organized as follows. Section 2 outlines existing PaI methods and describes how to calculate a compression ratio that correctly compares them with MiCA. Then, Section 3 proposes MiCA, and Section 4 compares MiCA with the RPI and RbPI methods. Finally, Section 5 concludes this paper.

2 Related Work

PaI algorithms can be categorized into two groups: 1) those that learn the criterion of pruning before training weights (i.e., RbPI); and 2) those that prune randomly (i.e., RPI). The methods included in these groups can be compared by using the compression ratio obtained by eliminating unused edges (Vysogorets & Kempe, 2023). This section recapitulates the method for calculating the corrected compression ratio and outlines the literature on RPI and RbPI.

Calculation of Corrected Compression Ratios. Recently, Vysogorets & Kempe (2023) found that PaI algorithms produce significant amounts of redundant parameters that can be removed without affecting the output. Figure 1 illustrates this phenomenon. A pruned network has 10 edges, but the 4 dashed edges do not affect the output. Thus, the apparent compression ratio is $21/10 = 2.1$, but it can also be regarded as $21/(10 - 4) = 3.5$. Correcting the compression ratio calculation by removing them results in a fairer comparison between subnetworks. Subsequent sections use this compression ratio calculation.

Random Pruning at Initialization (RPI). RPI prunes each layer randomly based on a pre-defined sparsity distribution calculated by a pre-defined compression ratio. To date, various sparsity distribution

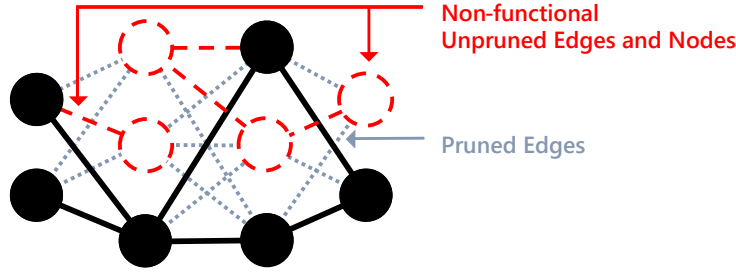


Figure 1: Some edges and nodes stop affecting the output (dashed line) as the network becomes sparser by pruning dotted lines.

design methods have been proposed. For example, Erdős-Rényi-Kernel (ERK) (Evci et al., 2020), which was devised in the context of random graphs, determines the density of l -th layer to be proportional to the scale $(C_{in}^{(l)} + C_{out}^{(l)} + k_h^{(l)} + k_w^{(l)}) / (C_{in}^{(l)} \times C_{out}^{(l)} \times k_h^{(l)} \times k_w^{(l)})$, where $C_{in}^{(l)}$, $C_{out}^{(l)}$, $k_h^{(l)}$, and $k_w^{(l)}$ denote input channels, output channels, kernel height, and kernel width of the l -th layer, respectively. Ideal Gas Quotas (IGQ) (Vysogorets & Kempe, 2023) focuses on the fact that traditional global pruning methods (Lee et al., 2019; 2021; Tanaka et al., 2020) intensively remove parameter-heavy layers. It determines the constant F based on the target compression ratio and calculates the density of l -th layer as $(F |E^{(l)}| + 1)^{-1}$, where $|E^{(l)}|$ is the number of edges in the l -th layer. The subnetworks to be pruned based on these sparsity distributions achieve comparable or better performance against RbPI and RPI using other distributions (Vysogorets & Kempe, 2023). However, those with skip connection-free architectures cannot achieve such performance as the compression ratio increases. Although some solutions, such as adding edges and resampling (Gadhikar et al., 2023), are proposed to address this RPI weakness, they are impossible at higher compression ratios or inefficient. Unlike these methods, our approach is efficient and works at higher compression ratios. It constructs a subnetwork with only pre-allocated edges, except that pre-defined sparsity distributions invariably cause non-functional neurons (i.e., nodes). Moreover, our approach does not require the iterative pruning operation.

Ranking-Based Pruning at Initialization (RbPI). RbPI determines the pruning priorities based on the initial state of a network and a dataset. For example, SNIP (Lee et al., 2019) uses the magnitude of the backpropagation gradient after one iteration as a parameter’s pruning priority. GraSP (Wang et al., 2020) prunes the edges that do not reduce the gradient flow of a subnetwork for a dataset preferentially. SynFlow (Tanaka et al., 2020) updates a parameter’s pruning priority by using the l_1 -path norm of a network (Neyshabur et al., 2015) as a loss and prunes iteratively without dataset. Unlike other RbPI methods, it can avoid layer-collapse (Hayou et al., 2020) at high compression ratios. As seen from these, RbPI takes into account the information flow of the initialized network for the pruning criteria. Hence, its sparse sub-networks tend to connect layers to each other. On the other hand, RbPI requires pre-training to calculate priorities, which is prohibitively expensive. RPI and our approach have a lower cost because they do not need to calculate priorities.

3 MiCA: Minimum Connection Assurance

Recently, Vysogorets & Kempe (2023) showed that RPI degrades accuracy significantly in skip connection-free architectures in high compression range. Gadhikar et al. (2023) focused on the connection of each layer to improve accuracy in the situation, but those methods cannot put higher compression ratios into perspective or are inefficient due to iterative operations. Given the importance of the connection of each layer as revealed by Gadhikar et al. (2023), this section introduces a novel PaI method, MiCA, which is completed in a single procedure and takes into account the case of extremely high compression ratios. MiCA constructs a subnetwork using the minimum connection as a frame. In order to create a minimum connection, it is essential to determine how many nodes are used in each layer based on the architecture and pre-allocated

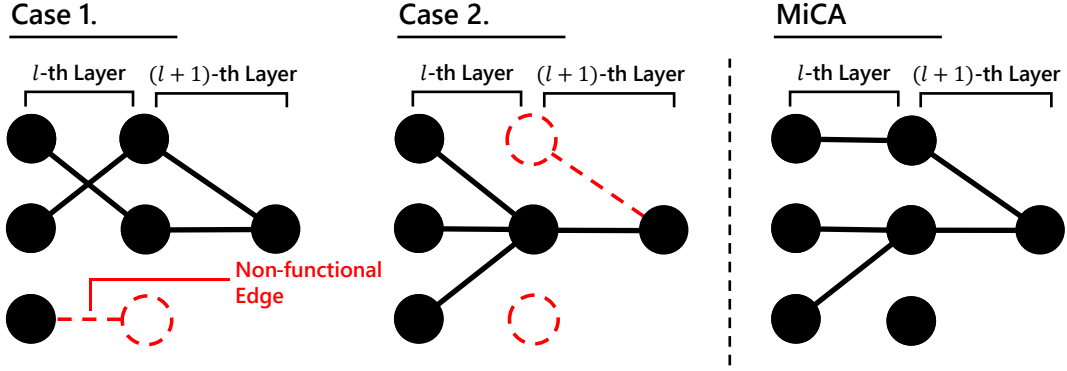


Figure 2: In cases 1 and 2, randomly pruning the network can make some edges non-functional. On the other hand, MiCA keeps all edges functional while allowing randomness in the connections.

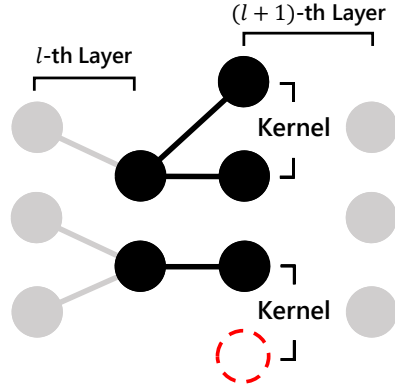


Figure 3: A node in l -th convolutional layer has a connection with up to kernel size $k_h^{(l)} \times k_w^{(l)}$ input nodes in $(l+1)$ -th convolutional layer. A convolution process makes these connections and does not use network parameters.

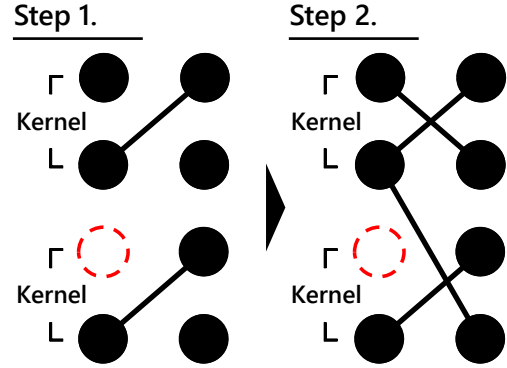


Figure 4: Edge placement procedures to construct a minimum connection. Each kernel is connected to an output node, and then unconnected nodes are connected.

edges. Therefore, we first introduce the way to analyze the number of nodes. Then, we describe how to create a minimum connection with a small amount of edges using the analysis result.

Figure 2 exemplifies two situations that are suboptimal for training or inference with a sparse network: in case 1, the edges of l -th layer connect to almost all nodes of $(l+1)$ -th layer. It causes some non-functional edges at l -th layer if the number of edges of the $(l+1)$ -th layer is low; in case 2, almost all edges connect to a few nodes of the next layer. It will cause some non-functional edges at $(l+1)$ -th layer if the number of edges of the $(l+1)$ -th layer is high. These can happen when pre-determining the number of edges in each layer and pruning a model randomly. We pre-determine which nodes to use to avoid these dire situations based on a pre-defined sparsity distribution and each layer architecture.

This section considers a convolutional neural network constructed with L layers. Let $|C_{in}^{(l)}|$, $|C_{out}^{(l)}|$, $k_h^{(l)}$, and $k_w^{(l)}$ be the number of input channels, the number of output channels, height, and width of the kernel of l -th layer, respectively. Note that these values are pre-determined in the architecture of the layer. Then, we define the input nodes and the output nodes in the pruned l -th layer as $V_{in}^{(l)}$ and $V_{out}^{(l)}$, respectively. Pruned l -th convolutional layer can be regarded as a bipartite graph formed by connecting $V_{in}^{(l)}$ and $V_{out}^{(l)}$. Thus, given l -th layer edges $E^{(l)}$, it is needed to satisfy $\max(|V_{in}^{(l)}|, |V_{out}^{(l)}|) \leq |E^{(l)}| \leq |V_{in}^{(l)}| |V_{out}^{(l)}|$ to use all edges.

Solving this inequality for $|V_{in}^{(l)}|$ leads to

$$\frac{|E^{(l)}|}{|V_{out}^{(l)}|} \leq |V_{in}^{(l)}| \leq |E^{(l)}|. \quad (1)$$

Also, as shown in Figure 3, up to $k_h^{(l)} \times k_w^{(l)}$ input nodes of the l -th layer connect to an output node of the $(l-1)$ -th layer. Hence, $|V_{in}^{(l)}|$ satisfies the inequality

$$|V_{out}^{(l-1)}| \leq |V_{in}^{(l)}| \leq |V_{out}^{(l-1)}| \times k_h^{(l)} \times k_w^{(l)}. \quad (2)$$

From these two inequalities, we can use all edges if $|V_{out}^{(l-1)}|$ fulfills

$$\frac{|E^{(l)}|}{|V_{out}^{(l)}|} \leq |V_{out}^{(l-1)}| \leq \frac{|E^{(l)}|}{k_h^{(l)} \times k_w^{(l)}}. \quad (3)$$

Then, we select $|V_{out}^{(l-1)}|$ randomly from the range. In contrast, the inequality does not hold if $|V_{out}^{(l)}| < k_h^{(l)} \times k_w^{(l)}$. In this case, we determine $|V_{out}^{(l-1)}|$ as $\left\lceil \frac{|E^{(l)}|}{|V_{out}^{(l)}|} \right\rceil$ to minimize the loss of input nodes $|V_{in}^{(l)}|$.

As shown in Equation 3, it is necessary to know $|V_{out}^{(l)}|$ for setting $|V_{out}^{(l-1)}|$. Thus, we sequentially analyze the number of connectable nodes from the output layer to the input layer. For the L -th layer (e.g., a classifier), the number of output nodes $|V_{out}^{(L)}|$ takes the number of classes. If the network has a branching structure, such as residual connections, we select the larger $|V_{out}^{(l-1)}|$ obtained at each branch to keep all input nodes connected. The computational cost of this analysis is negligibly tiny because it is completed using only simple operations.

Once the nodes used in each layer are determined, the subnetwork construction of MiCA is completed by placing the edges. First, we construct a minimum connection in two steps, as shown in Figure 4:

1. Select an input node from a kernel and connect it to an unconnected output node. Note that it is performed for all kernels.
2. Connect the unconnected nodes. If the minimum degree of input or output nodes is already 1, we connect the unconnected nodes to randomly selected nodes.

Step 1 is particularly essential when the compression ratio is extremely high (i.e., there are a few edges to place). For a network between convolutional layers, as shown in Figure 3, using one edge for each kernel first prevents the edges in the previous layer from becoming non-functional. These steps are completed sequentially from the input layer to the output layer as we connect the nodes used in the input side to the output node. Finally, we place the remaining edges randomly within the nodes in the minimum connection. These edges work wherever they are placed because minimum connection exists.

4 Experiments and Results

This section evaluates MiCA on image classification using the experimental setup described in Section 4.2. It shows that MiCA performs better than conventional RPI and RbPI methods, especially in the high compression range and regardless of whether skip connections exist.

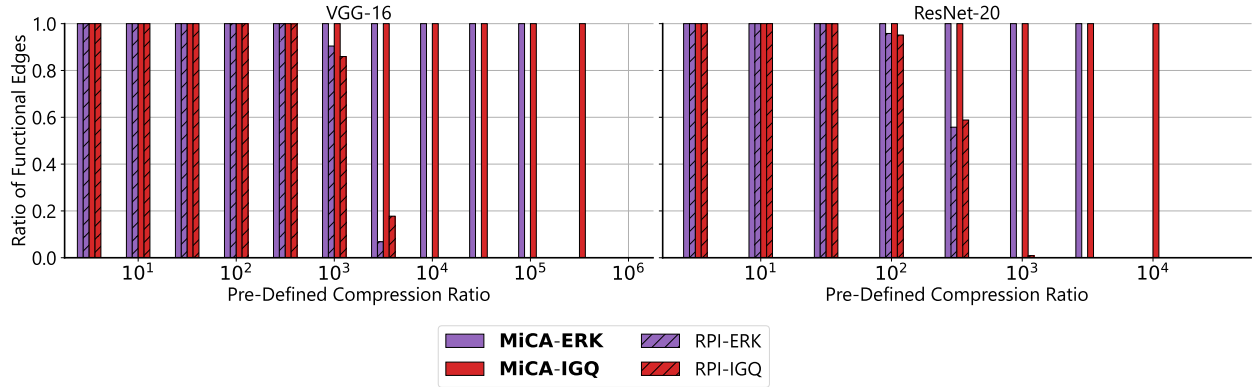


Figure 5: Comparison of the ratio of functional edges between RPI and MiCA with ERK and IGQ for pre-defined compression ratios. Unlike other RPI methods, MiCA makes all edges functional, even at a high compression ratio.

4.1 Notation of Methods

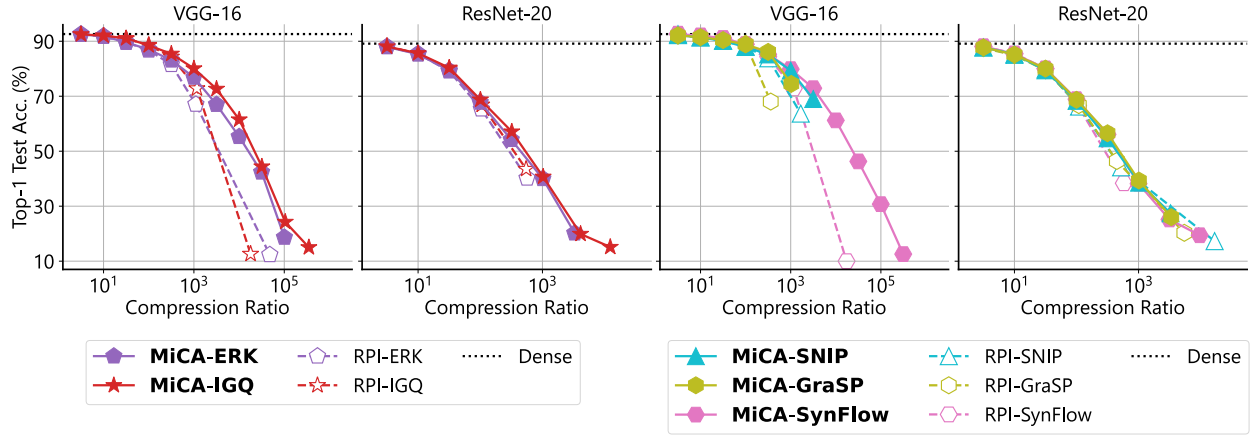
We briefly introduce the notation of the methods compared in the subsequent experiments in advance. RPI uses a pre-defined sparsity distribution for ERK, IGQ, SNIP, GraSP, and SynFlow. Here, we consider the sparsity distribution of the network pruned by SNIP, GraSP, and SynFlow, which are RbP, as a pre-defined sparsity distribution. RPI methods with these distributions are denoted **RPI-ERK**, **RPI-IGQ**, **RPI-SNIP**, **RPI-GraSP**, and **RPI-SynFlow**, respectively. Similarly, MiCA using these distributions are denoted **MiCA-ERK**, **MiCA-IGQ**, **MiCA-SNIP**, **MiCA-GraSP**, and **MiCA-SynFlow**, respectively. For the RbPI experiments, SNIP, GraSP, and SynFlow are chosen as RbPI methods, and these are specified as **RbPI-SNIP**, **RbPI-GraSP**, **RbPI-SynFlow**, respectively. For more detail on each RPI and RbPI method, see Section 2.

4.2 Experimental Settings

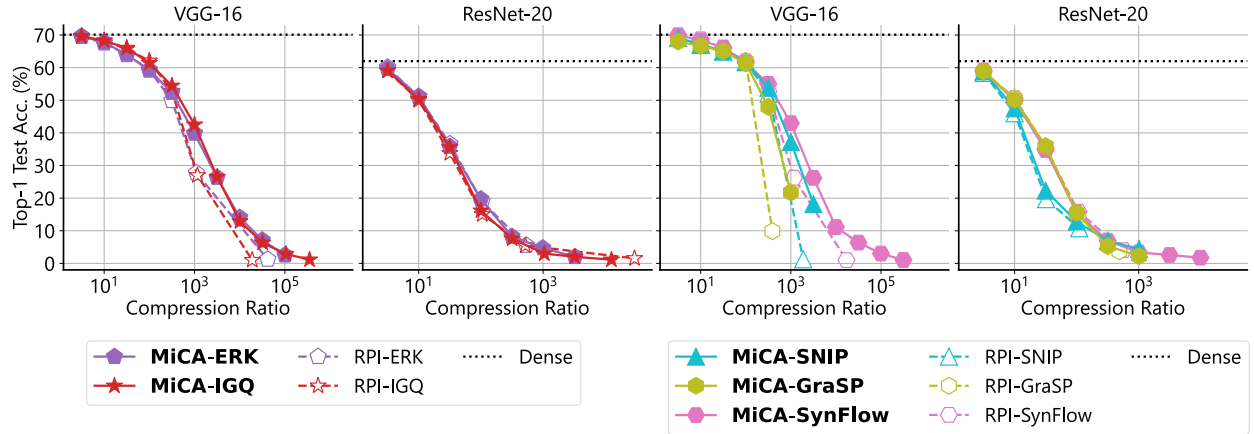
This paper evaluates MiCA in image classification on the CIFAR-10 and CIFAR-100 datasets. For CIFAR-10 and CIFAR-100, 40,000 images are used as training data and 10,000 as validation data. The architectures used in the experiments are VGG-16 and ResNet-20. Implementation is based on the code provided by Tanaka et al. (2020). In particular, VGG-16 includes a batch normalization layer and removes the bias of the convolutional layer. All experiments use stochastic gradient descent (SGD) applying Nesterov’s acceleration method (Nesterov, 1983) with a momentum of 0.9. We use values $10^{0.5}, 10^1, \dots, 10^{5.5}$, and 10^6 as a pre-defined compression ratio. Note that results are not plotted if all edges are non-functional after compression. CIFAR-10 and CIFAR-100 experiments are run three times with a batch size 128 for 160 epochs. For VGG-16, the weight decay is set to 0.0001, and the learning rate is started at 0.1 and multiplied by 0.1 after 60 and 120 epochs. For ResNet-20, the weight decay is set to 0.0005, and the learning rate is started at 0.01 and multiplied by 0.2 after 60 and 120 epochs. SNIP and GraSP use $10\times$ the training data amount relative to the dataset’s number of classes and a batch size of 128. SynFlow prunes the initialized network for 100 iterations.

4.3 Pre-Defined Sparsity Distribution Maintenance of MiCA

This section shows how much MiCA and RPI leave functional edges for each pre-defined compression ratio. The result is plotted in Figure 5. The ratios of functional edges for RPI-ERK and RPI-IGQ begin to decrease around $10^3\times$ compression ratio for VGG-16 and $10^2\times$ for ResNet-20. On the other hand, MiCA continues to use all the pre-allocated edges as much as possible and maintains the sparsity distribution even in the high compression range. MiCA-ERK and MiCA-IGQ have almost all functional edges even at $10^5\times$ compression ratio for VGG-16 and $10^{3.5}\times$ compression ratio for ResNet-20. However, all edges can be non-functional at high compression ratios, as seen in MiCA-ERK for VGG-16 at $10^{5.5}\times$ compression ratio. This phenomenon



(a) CIFAR-10 experiments.



(b) CIFAR-100 experiments.

Figure 6: Comparison of accuracy between RPI and MiCA on CIFAR-10 and CIFAR-100. The two columns on the left compare RPI and MiCA with ERK and IGQ, and the two columns on the right compare RPI and MiCA using the sparsity distribution of SNIP, GraSP, and SynFlow. In VGG-16, MiCA reduces performance degradation in the high compression range and improves the accuracy and compression ratio trade-offs. It also shows slight performance improvement for the sparsity distributions such as IGQ and SynFlow in ResNet-20 experiments using CIFAR-10. Note that some plots are not plotted because the corrected compression ratio is infinite.

is caused by the way the pre-defined sparsity distribution is designed. Some pre-defined sparsity distributions allocate no edges to a few layers when the pre-defined compression ratio is exceptionally high. As a result, all edges are non-functional, regardless of how they are placed.

4.4 MiCA vs. Random Pruning

Figure 6 compares MiCA and RPI with ERK, IGQ, SNIP, GraSP, and SynFlow. First, we focus on the ERK and IGQ experiments of the two left columns. In the VGG-16 experiment (first column), RPI-ERK and RPI-IGQ have a sharp performance drop for $\geq 10^3 \times$ compression ratios for both CIFAR-10 and CIFAR-100. These with a pre-defined compression ratio of $10^{3.5} \times$ have an actual compression ratio of $> 10^4 \times$, and those with a pre-defined compression ratio of $10^4 \times, 10^{4.5} \times, \dots, 10^5 \times$, and $10^6 \times$ cannot be plotted as all edges are non-functional (i.e., the corrected compression ratio is infinite). On the other hand, MiCA-ERK and MiCA-

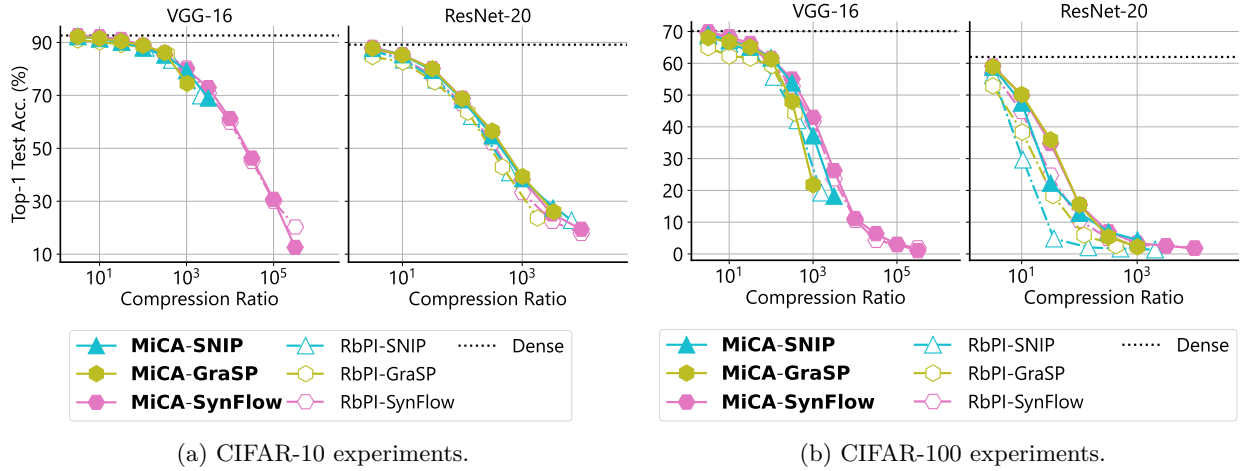


Figure 7: Comparison of accuracy between RbPI and MiCA. MiCA improves the accuracy and compression ratio trade-off more than RbPI in ResNet-20. MiCA also performs as well as RbPI in the VGG-16 experiments. Interestingly, MiCA outperforms RbPI in both VGG-16 and ResNet-20 experiments using CIFAR-100.

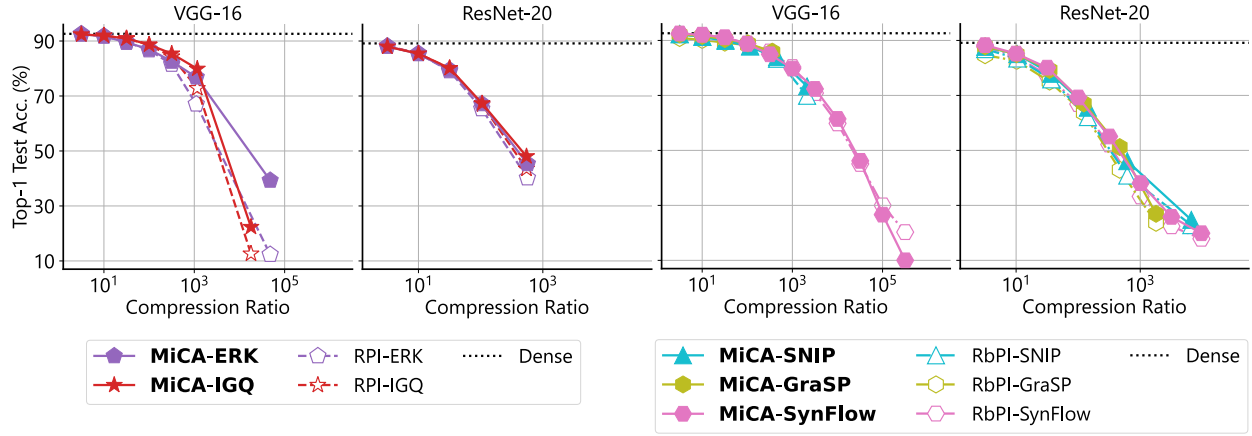
IGQ maintain the same compression ratio as the pre-defined compression ratio and suffer less performance degradation. Note that MiCA-ERK is not plotted for $10^{5.5} \times$ and $10^6 \times$ compression ratios. This phenomenon is due to ERK’s design, as mentioned in Section 4.3, and the same occurs in other sparsity distributions (e.g., IGQ) and other experiments. As shown in Figure 6 (a), MiCA-IGQ achieves an accuracy of 44.4% for $10^{4.5} \times$ compression ratio, significantly higher than RPI-IGQ’s accuracy of 12.6% for $10^{4.3} \times$ compression ratio. This result suggests that the minimum connection supports learning in the high compression range.

In the ResNet-20 experiment (second column), the performance difference between MiCA and RPI is less drastic than in the VGG-16 experiment. In particular, CIFAR-100 experiments (Figure 6 (b)) show little difference. However, it is hardly surprising considering that skip connections help randomly pruned networks learn in the high compression range (Hoang et al., 2023). For CIFAR-10 experiments (Figure 6 (a)), MiCA-ERK slightly improves the trade-off between compression ratio and accuracy against RPI-ERK in the compression range of 10^2 – $10^3 \times$. ResNet-20 has several layers that do not have skip connections, and it is therefore considered that the minimum connection supports learning in those layers.

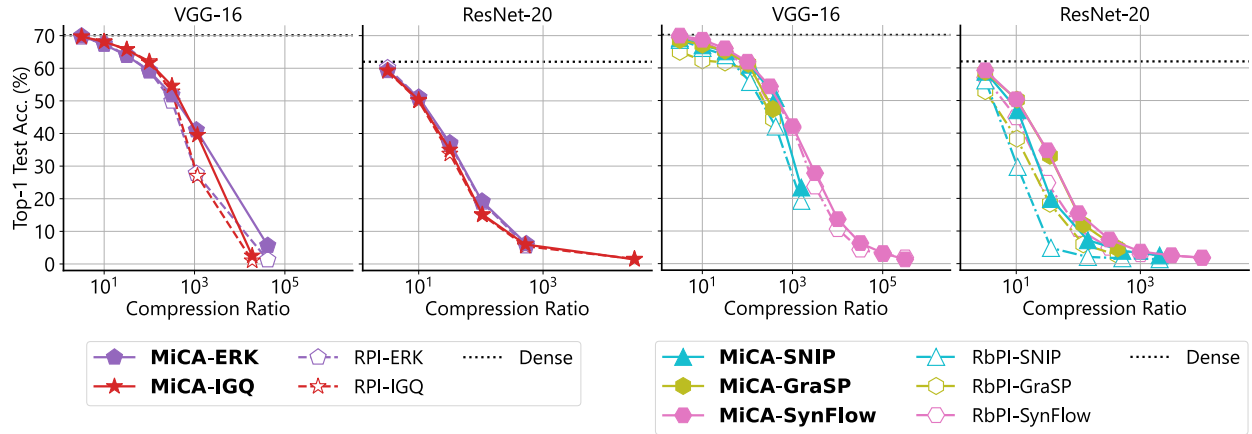
Then, we state the results of the right two columns in Figure 6. In the VGG-16 experiment (third column), MiCA-SNIP, MiCA-GraSP, and MiCA-SynFlow significantly improve the trade-off between accuracy and compression ratio as in the ERK and IGQ experiments. It is particularly evident in the CIFAR-100 experiment (Figure 6 (b)). Thus, MiCA overcomes the performance degradation at a high compression range for skip connection-free architectures, regardless of the pre-defined sparsity distribution. However, the performance improvements in the ResNet-20 (fourth column) are still minute. This result seems regrettable but highlights the importance of skip connections to RPI.

4.5 MiCA vs. Ranking-Based Pruning

Figure 7 compares the compression ratio and accuracy between RbPI and MiCA on CIFAR-10 (left two columns) and CIFAR-100 (right two columns). We employ SNIP, GraSP, and SynFlow as RbPI methods. Despite random connections, MiCA achieves comparable performance to RbPI in VGG-16. Although it has already been observed that randomly pruned networks achieve performance comparable to RbPI (Frankle et al., 2021), to the best of our knowledge, this is the first time that the same result is reported in the high compression range. In other words, this result shows that RbPI only learns layer-to-layer connections, not high-performance subnetworks, regardless of the compression ratio. Interestingly, the ResNet-20 experiments show a more pronounced performance difference than the RPI result in Figure 6. It suggests that MiCA’s



(a) CIFAR-10 experiments.



(b) CIFAR-100 experiments.

Figure 8: The accuracy of RPI and MiCA with the same sparsity distribution is compared on CIFAR-10 and CIFAR-100. When the number of parameters used in each layer is matched, MiCA performs better than other PaI methods when the compression ratio is high.

RPI aspect helps it to maintain higher performance than RbPI in ResNet-20 because MiCA connects each layer randomly.

4.6 Performance Comparison for the Same Sparsity Distribution

Previous sections match the pre-defined sparsity distribution of MiCA and other PaI methods and compare each method. In contrast, this section evaluates the corrected sparsity distribution of RPI and RbPI as the pre-defined sparsity distribution of MiCA. Then, we show that the network structure constructed by MiCA achieves higher performance against other PaI methods, even when the number of edges at each layer is matched.

Figure 8 compares RPI, RbPI, and MiCA for the same sparsity distribution. The two columns on the left show RPI experiments, while the two on the right show RbPI experiments. Even when using the same sparsity distribution, MiCA improves the accuracy more than RPI (left two columns). The results in VGG-16 on CIFAR-10 are remarkable: MiCA-ERK achieves 27.0% higher accuracy than RPI-ERK with $10^{4.7} \times$ compression ratio. Performance improvements of MiCA can also be seen in ResNet-20 at a high compression ratio. Both RPI and MiCA ought to make similar networks due to the random connection in this situation,

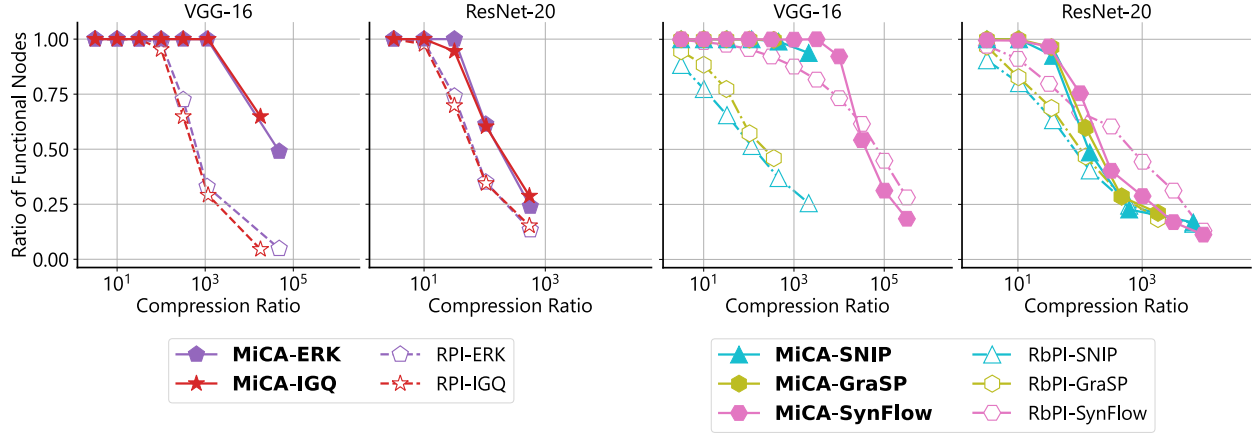


Figure 9: Comparison of the compression ratio with the ratio of functional nodes among the nodes flowing information from the first layer. MiCA maintains a high functional node ratio even in a high compression range compared to other PaI methods. Exceptions are RbPI-SynFlow, which retains a higher ratio than MiCA at high compression ratios.

but the performance difference is more noticeable when the compression ratio is high. In addition, the subnetworks built by RbPI suffer more performance degradation than those built by MiCA, regardless of compression ratio and architecture (right two columns). However, in the VGG-16 experiments (third column), RbPI-SynFlow is more accurate than MiCA in the compression range above $10^5 \times$. At $10^{5.5} \times$ compression ratio in CIFAR-10 experiments, MiCA-SynFlow achieves 10.0%, which is not different from the random performance, while RbPI-SynFlow achieves 20.3%. Furthermore, these accuracies are almost identical for CIFAR-100 experiments. In other words, RbPI-SynFlow is superior to MiCA only when the compression ratio is extremely high and the architecture does not have skip connections.

4.7 Top-to-Bottom Information Propagation

This section shows that MiCA propagates information without loss compared to other PaI methods. As in the previous section, the pre-defined sparsity distribution of MiCA is matched to the corrected sparsity distribution of each method.

Figure 9 compares the ratio of functional nodes to all nodes that information flows from the first layer for each compression ratio. While RPI and RbPI methods reduce the ratio of functional nodes as the compression ratio increases, MiCA maintains the high ratio of functional nodes. In particular, RPI-ERK and RPI-IGQ have almost 3/4 of the nodes non-functional in VGG-16 at $10^3 \times$ compression ratio, whereas MiCA-ERK and MiCA-IGQ keep all nodes functional. Interestingly, RbPI-SynFlow keeps the ratio of functional nodes relatively higher than other PaI methods without MiCA, even in the high compression range. In contrast, RbPI methods make nodes non-functional even in the low compression range. At the compression ratio of $> 10^4 \times$ in VGG-16 and $> 10^2 \times$ in ResNet-20, RbPI has a higher ratio of functional nodes than MiCA. Given that RbPI-SynFlow achieves comparable or higher accuracy than MiCA at $\geq 10^5 \times$ compression ratio in VGG-16 (Figure 8) and that RbPI builds a specific network at the skip connection-free architecture (Hoang et al., 2023), it suggests that sparse networks with a high ratio of functional nodes maintain accuracy at a high compression ratio. MiCA makes nodes non-functional in a high compression range, but the ratio is lower than in other methods. At $10^{4.7} \times$ compression ratio in VGG-16, MiCA-ERK keeps nearly 50% of the nodes functional, whereas RPI-ERK keeps most of the nodes non-functional. In the ResNet-20 experiments, MiCA also has a higher ratio of functional nodes than other methods, but the difference is lower than in VGG-16. It shows that information from the first layer flows to the subsequent layers even after pruning due to skip connections. In other words, the top-to-bottom information flow is narrowed by pruning, but skip connections allow it to flow to the subsequent layers.

Why do some nodes become non-functional even after recalculating the compression ratio? This is because the convolution process produces non-functional nodes. The convolution process connects one input node in l -th convolutional layer with several output nodes in $(l + 1)$ -th convolutional layer, as shown in Figure 3. This connection does not use network parameters; hence, it is preserved after recalculating the compression ratio. Consequently, nodes can be non-functional if the number of remaining edges in $(l + 1)$ -th layer is small.

5 Conclusion and Future Work

This paper proposes Minimum Connection Assurance (MiCA), a restricted random pruning at initialization method. MiCA combines a departure from costly connection learning by random connections and the prevention of rapid accuracy degradation in the high compression range by preserving connections between layers. In other words, it takes the best of both RPI and RbPI’s respective strengths. It improves the accuracy and compression ratio trade-off of conventional PaI methods on the VGG and ResNet architectures. In particular, MiCA allows training at high compression ratios, which RPI could not achieve on the architectures without skip connections. Our experiments find that random connections can substitute the connections learned by RbPI. Although this has been demonstrated previously in the low compression range, this work is the first to demonstrate it in the high compression range. The analysis of the experimental results shows that MiCA sends information flowing from the first layer to the subsequent layers without loss compared to other PaI methods. It suggests that controlling the spread of information at each layer and sending information to the bottom layers without loss by pruning improves accuracy, even if the pruning is random.

A future direction of this work is to investigate sparsity distributions regarding minimum connection. Although MiCA achieves higher performance than conventional PaI methods, the accuracy in terms of compression ratio depends on the assigned sparsity distribution. Therefore, it is possible to discover a novel sparsity distribution by examining the importance of each layer, assuming a condition where the connection is maintained.

References

- Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *Proc. Int. Conf. Mach. Learn.*, pp. 322–332. PMLR, 2019.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Proc. Adv. Neural Inform. Process. Syst.*, 33:1877–1901, 2020.
- Yves Chauvin. A back-propagation algorithm with optimal use of hidden units. *Proc. Adv. Neural Inform. Process. Syst.*, 1, 1988.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *Proc. Int. Conf. Learn. Repr.*, 2021.
- Utku Evci, Trevor Gale, Jacob Menick, Pablo Samuel Castro, and Erich Elsen. Rigging the lottery: Making all tickets winners. In *Proc. Int. Conf. Mach. Learn.*, pp. 2943–2952. PMLR, 2020.
- Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *Proc. Int. Conf. Learn. Repr.*, 2019.
- Jonathan Frankle, Gintare Karolina Dziugaite, Daniel M Roy, and Michael Carbin. Stabilizing the lottery ticket hypothesis. *arXiv preprint arXiv:1903.01611*, 2019.
- Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. Pruning neural networks at initialization: Why are we missing the mark? In *Proc. Int. Conf. Learn. Repr.*, 2021.

- Advait Harshal Gadhikar, Sohom Mukherjee, and Rebekka Burkholz. Why random pruning is all we need to start sparse. In *Proc. Int. Conf. Mach. Learn.*, pp. 10542–10570. PMLR, 2023.
- Babak Hassibi, David G Stork, and Gregory J Wolff. Optimal brain surgeon and general network pruning. In *IEEE international conference on neural networks*, pp. 293–299. IEEE, 1993.
- Soufiane Hayou, Jean-Francois Ton, Arnaud Doucet, and Yee Whye Teh. Pruning untrained neural networks: Principles and analysis. *arXiv preprint arXiv:2002.08797*, 2020.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. IEEE Comput. Soc. Conf. Comput. Vis. and Pattern Recognit.*, pp. 770–778, 2016.
- Duc N.M Hoang, Shiwei Liu, Radu Marculescu, and Zhangyang Wang. Revisiting pruning at initialization through the lens of ramanujan graph. In *Proc. Int. Conf. Learn. Repr.*, 2023.
- Masumi Ishikawa. Structural learning with forgetting. *Neural networks*, 9(3):509–521, 1996.
- Siddhant Jayakumar, Razvan Pascanu, Jack Rae, Simon Osindero, and Erich Elsen. Top-kast: Top-k always sparse training. *Proc. Adv. Neural Inform. Process. Syst.*, 33:20744–20754, 2020.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Yann LeCun, John Denker, and Sara Solla. Optimal brain damage. *Proc. Adv. Neural Inform. Process. Syst.*, 2, 1989.
- Jaeho Lee, Sejun Park, Sangwoo Mo, Sungsoo Ahn, and Jinwoo Shin. Layer-adaptive sparsity for the magnitude-based pruning. In *Proc. Int. Conf. Learn. Repr.*, 2021.
- Namhoon Lee, Thalaiyasingam Ajanthan, and Philip Torr. SNIP: Single-shot network pruning based on connection sensitivity. In *Proc. Int. Conf. Learn. Repr.*, 2019.
- Shiwei Liu, Tianlong Chen, Xiaohan Chen, Li Shen, Decebal Constantin Mocanu, Zhangyang Wang, and Mykola Pechenizkiy. The unreasonable effectiveness of random pruning: Return of the most naive baseline for sparse training. In *Proc. Int. Conf. Learn. Repr.*, 2022.
- Decebal Constantin Mocanu, Elena Mocanu, Peter Stone, Phuong H Nguyen, Madeleine Gibescu, and Antonio Liotta. Scalable training of artificial neural networks with adaptive sparse connectivity inspired by network science. *Nature communications*, 9(1):2383, 2018.
- Yurii Evgen’evich Nesterov. A method of solving a convex programming problem with convergence rate $O(1/\sqrt{k})$. In *Doklady Akademii Nauk*, volume 269, pp. 543–547. Russian Academy of Sciences, 1983.
- Behnam Neyshabur, Russ R Salakhutdinov, and Nati Srebro. Path-SGD: Path-normalized optimization in deep neural networks. *Proc. Adv. Neural Inform. Process. Syst.*, 28, 2015.
- Behnam Neyshabur, Zhiyuan Li, Srinadh Bhojanapalli, Yann LeCun, and Nathan Srebro. The role of over-parametrization in generalization of neural networks. In *Proc. Int. Conf. Learn. Repr.*, 2019.
- Ilan Price and Jared Tanner. Dense for the price of sparse: Improved performance of sparsely initialized networks via a subspace offset. In *Proc. Int. Conf. Mach. Learn.*, pp. 8620–8629. PMLR, 2021.
- Alex Renda, Jonathan Frankle, and Michael Carbin. Comparing rewinding and fine-tuning in neural network pruning. In *Proc. Int. Conf. Learn. Repr.*, 2020.
- Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-LM: Training multi-billion parameter language models using model parallelism. *arXiv preprint arXiv:1909.08053*, 2019.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

- Jingtong Su, Yihang Chen, Tianle Cai, Tianhao Wu, Ruiqi Gao, Liwei Wang, and Jason D Lee. Sanity-checking pruning methods: Random tickets can win the jackpot. *Proc. Adv. Neural Inform. Process. Syst.*, 33:20390–20401, 2020.
- Hideori Tanaka, Daniel Kunin, Daniel L Yamins, and Surya Ganguli. Pruning neural networks without any data by iteratively conserving synaptic flow. *Proc. Adv. Neural Inform. Process. Syst.*, 33:6377–6389, 2020.
- Artem Vysogorets and Julia Kempe. Connectivity matters: Neural network pruning through the lens of effective sparsity. *J. Mach. Learn. Res.*, 24:99–1, 2023.
- Chaoqi Wang, Guodong Zhang, and Roger Grosse. Picking winning tickets before training by preserving gradient flow. In *Proc. Int. Conf. Learn. Repr.*, 2020.
- Andreas Weigend, David Rumelhart, and Bernardo Huberman. Generalization by weight-elimination with application to forecasting. *Proc. Adv. Neural Inform. Process. Syst.*, 3, 1990.
- Kaiyue Wen, Jiaye Teng, and Jingzhao Zhang. Benign overfitting in classification: Provably counter label noise with larger models. In *Proc. Int. Conf. Learn. Repr.*, 2022.
- Sanghyun Woo, Shoubhik Debnath, Ronghang Hu, Xinlei Chen, Zhuang Liu, In So Kweon, and Saining Xie. Convnext v2: Co-designing and scaling convnets with masked autoencoders. In *Proc. IEEE Comput. Soc. Conf. Comput. Vis. and Pattern Recognit.*, pp. 16133–16142, 2023.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.
- Guodong Zhang, James Martens, and Roger B Grosse. Fast convergence of natural gradient descent for over-parameterized neural networks. *Proc. Adv. Neural Inform. Process. Syst.*, 32, 2019.