

Effective Sparsification of Neural Networks with Global Sparsity Constraint

Xiao Zhou^{1*}, Weizhong Zhang^{1*}, Hang Xu², Tong Zhang¹

¹The Hong Kong University of Science and Technology ²Huawei Noah's Ark Lab

xzhoubi@connect.ust.hk, weizhong@ust.hk, xu.hang@huawei.com, tongzhang@tongzhang-ml.org

Abstract

Weight pruning is an effective technique to reduce the model size and inference time for deep neural networks in real-world deployments. However, since magnitudes and relative importance of weights are very different for different layers of a neural network, existing methods rely on either manual tuning or handcrafted heuristic rules to find appropriate pruning rates individually for each layer. This approach generally leads to suboptimal performance. In this paper, by directly working on the probability space, we propose an effective network sparsification method called probabilistic masking (ProbMask), which solves a natural sparsification formulation under global sparsity constraint. The key idea is to use probability as a global criterion for all layers to measure the weight importance. An appealing feature of ProbMask is that the amounts of weight redundancy can be learned automatically via our constraint and thus we avoid the problem of tuning pruning rates individually for different layers in a network. Extensive experimental results on CIFAR-10/100 and ImageNet demonstrate that our method is highly effective, and can outperform previous state-of-the-art methods by a significant margin, especially in the high pruning rate situation. Notably, the gap of Top-1 accuracy between our ProbMask and existing methods can be up to 10%. As a by-product, we show ProbMask is also highly effective in identifying supermasks, which are sub-networks with high performance in a randomly weighted dense neural network.

1. Introduction

Weight pruning [9] is a popular technique for alleviating the weight redundancy in deep neural networks (DNNs) to improve inference efficiency and decrease computation demands. Typical pruning algorithms usually prune the unimportant weights by developing proper criteria. It is repeatedly reported in the literature [8, 22, 39, 21] that by pruning one can reduce the neural network size and improve

the inference efficiency significantly with quite slight or even negligible loss on performance, which makes deploying large-scale DNNs on equipment with limited computational and memory budget possible.

What can serve as a suitable global comparator to measure weight importance and identify sparsity level for different layers is a long-standing problem [7] though impressive results have been achieved. We know that the core module in pruning is the explicit or implicit criterion for identifying the redundant weights, and it is difficult to develop a global criterion for the weights in all the layers. For example, in [9], the authors propose a simple yet effective criterion, i.e., for each layer it prunes all the weights below a certain threshold in a fully trained network. The threshold is obtained by sorting weight by its magnitude and retrieving the weight magnitude at the target pruning rate. The criterion is weight magnitude in this case. Notice that the magnitudes of the weights across layers could be quite different and different layers could have different amount of redundancy. If we use a global threshold for all the layers, then almost all the weights in certain layers could be pruned in order to achieve high enough pruning ratio, which will be verified in Section 5. Thus, we need to set a proper threshold or pruning ratio for each layer individually. In the networks with numerous layers, it is very difficult and even impossible to find the optimal thresholds or pruning ratios for all the layers manually. One reasonable compromise for such dilemma is to set sparsity level uniformly for different layers. However, this results in imperfect weight allocation obviously and gives unsatisfactory results on high pruning rates.

In this paper, to address the above limitations, we propose an effective network sparsification method called *probabilistic masking* (ProbMask). Firstly, we know that network pruning can be naturally formulated into a problem of finding a sparse binary mask \mathbf{m} as well as the weights at the same time to minimize the empirical loss (1). If the component m_i is equal to 0, it means that the corresponding weight is pruned. However, it is a discrete optimization problem and hard to solve. We notice that if we view the components m_i in the mask as independent Bernoulli ran-

*Equal contribution

dom variables with probability s_i being 1 and probability $(1 - s_i)$ being 0 and reparameterize them w.r.t. its probability, then the loss in problem (1) would become continuous over the probability space. Due to the nature of probability, probability can be used as a global criterion in all the layers. Therefore, we can control the model size via forcing the sum of all the probabilities s_i of the mask smaller than a proper value, leading to a global sparsity constraint in the probability space. In this way, the discrete optimization problem (1) is transformed into a constrained expected loss minimization problem (2) over a probability space, which is continuous. Finally, we adopt the Gumbel-Softmax trick to solve the continuous problem. As the optimizer goes on, the probabilities s_i would converge to either 0 or 1, i.e., \mathbf{m} would become close to a deterministic sparse mask. Thus, a fully trained mask would have quite low variance, making the loss of the sampled sparse network according to \mathbf{m} close to the expected loss in problem (2). Another appealing feature of our proposed method is that the amount of weight redundancy in each layer can be identified automatically by our global sparsity constraint and thus we do not need to choose different pruning ratios for different layers.

Experimental results on network pruning and supermask [40] finding demonstrate that our method is much more effective than the state-of-the-art methods on both small scale datasets and large scale datasets and can outperform them with a significant margin when the pruning rate is high.

The contribution and novelty of ProbMask can be summarized as follows:

1) We provide evidence showing that probability can serve as a suitable global comparator to measure weight importance and identify sparsity level for different layers, which is a long-standing problem [7].

2) We present a natural formulation of global sparsity constraint, and an optimization method that is practically effective. Our solution fixes the training and testing performance discrepancy problem observed in practice, which led to the failure of previous methods [23] on ImageNet [7].

3) We demonstrate the effectiveness of using probability as global comparator on small-scale and large-scale problems and various models and achieve state-of-the-art results on Top-1 accuracy and accuracy-versus-FLOPS curve.

4) We show ProbMask can also serve as a powerful tool for identifying supermasks, which are subnetworks with high performance in a randomly weighted dense neural network, and we achieve state-of-the-art results on Top-1 accuracy on CIFAR-100 under high pruning rates.

Notations: Let $\|\cdot\|_0$, $\|\cdot\|_1$ and $\|\cdot\|_2$ be the ℓ_0 , ℓ_1 and ℓ_2 norm of a real valued vector, respectively. We denote $\mathbf{1} \in \mathbb{R}^n$ to be a vector with all components equal to 1. In addition, $\{0, 1\}^n$ is the set of n -dimensional vectors with each coordinate valued in $\{0, 1\}$.

2. Related Work

Below, we first review the related work on network pruning. Next we review training methods for obtaining sparse networks which can be divided into two groups: dense-to-sparse training and sparse-to-sparse training. Then we review some probability-based methods for obtaining sparse networks and point out some limitations to differentiate them from our work. Finally we review another line of research on Lottery Tickets Hypothesis, SuperMask and Foresight Pruning.

2.1. Network Pruning

Network Pruning [10, 8, 39, 21, 24, 14, 41, 17, 34, 30, 38] has been extensively studied in recent years to reduce the model size and improve the inference efficiency of deep neural networks. Since it is a widely-recognized property that modern neural networks are always over-parameterized, pruning methods are developed to remove unimportant parameters in the fully trained dense networks to alleviate such redundancy. According to the granularity of pruning, existing pruning methods can be roughly divided into two categories, i.e., unstructured pruning and structured pruning. The former one is also called weight pruning, which removes the unimportant parameters in an unstructured way, that is, any element in the weight tensor could be removed. The latter one removes all the weights in a certain group together, such as kernel and filter. Since structure is taken into account in pruning, the pruned networks obtained by structured pruning are available for efficient inference on standard computation devices. In both structured and unstructured pruning methods, their key idea is to propose a proper implicit or explicit criterion (e.g., magnitude of the weight [9, 10, 8, 41, 6, 28, 1, 26, 35], scores based on Hessian, momentum or gradient [5, 20, 39, 19, 11, 4]) to evaluate the importance of the weight, kernel or filter and then remove the unimportant ones. The results in the literature [8, 22, 39, 21, 30, 18, 5, 35] demonstrate that pruning methods can significantly improve the inference efficiency of DNNs with minimal performance degradation, making the deployment of modern neural networks on resource-limited devices possible.

2.2. Dense-to-sparse and Sparse-to-sparse Training

We follow the convention of [18] to divide training algorithms for obtaining sparse networks into two groups: dense-to-sparse training and sparse-to-sparse training. Dense-to-sparse training starts with a dense network and obtains a sparse network at the end of the training [10, 41, 27, 6, 30, 36, 32, 23, 35]. ProbMask belongs to the group of dense-to-sparse training. [9, 41, 6, 30] follows the idea of using weight magnitude as the criterion. [41] manually set a uniform sparsity budget for different layers. [30] achieves strong results but needs multiple rounds of pruning

and retraining. [36] assigns auxiliary scores to weights and use it as the criterion. [36] suffers from the bias induced by the approximation of the step function and will have gradient vanishing problem when using ReLU and SoftPlus as the approximator. This makes the auxiliary scores hard to act as a global criterion. [32, 23, 27] base its criterion on reparameterization of probability and have the most connections with our work. We will fully discuss them in the next subsection.

Sparse-to-sparse training starts with a sparse network and maintain the sparsity during training [1, 26, 28, 5, 4]. It uses criterion like weight magnitude, weight gradient magnitude, momentum of weight to reallocate sparsity through training. Conceptually, sparse-to-sparse training can reduce the computational cost during training but it is hard to take effect without the support of sparse convolution framework on GPU. The performance of sparse-to-sparse training generally falls behind dense-to-sparse training under the same setting as shown in [33, 18].

2.3. Probability-based Methods

Compared to directly treating probability as the trainable variable, [32, 23, 27] consider probability as the hidden state and optimize on another space by reparameterization of probability. [32, 23, 27] achieve strong empirical results while we observe several shortcomings in the reparameterization process. [32] approximates the gradient through sampling process by biased STE (straight-through estimator) [2] and represent probability as the output of a hard-sigmoid function which induces gradient vanishing. [23] reparameterizes w.r.t hard-concrete and [23] is reported to fail to work on ImageNet dataset because of the performance gap between training and testing phases [7]. [27] also exhibits gradient vanishing problem due to function pattern of KL divergence and performance gap between training and testing phases by generating the test model by a cut-off manner rather than sampling. We solve the aforementioned problems by directly optimizing over probability space. We empirically demonstrate that probabilities finally converge to either 0 or 1 after training in Section 5, leading to a binary mask and fixing the training and testing performance discrepancy problem observed in practice. Besides, by explicitly control the model size via global sparsity constraint, users don't need to tune regularization parameters to achieve desired model size, which is a missing feature in [32, 23, 27].

2.4. Lottery Tickets Hypothesis, Supermask and Foresight Pruning

Lottery Ticket Hypothesis was proposed in [6], which conjectures and verifies that there exists sparse subnetworks which can be trained directly to achieve even better performance than dense counterparts with less training time.

[40] further analyzes the conditions for such phenomenon to hold and propose supermask, which conveys an intriguing idea that a good mask is enough to achieve surprisingly good performance with randomly initialized weights. [29] further asks the question what's hidden in a randomly weighted neural network and proposes a more effective algorithm on finding supermasks. Lottery Ticket Hypothesis also sheds light on whether we could find such subnetwork without training a dense network. [20] propose the first algorithm to find such subnetworks and [33] further improves performance on high pruning rates. However, both of them could not achieve better performance than the latest pruning methods.

3. Effective Sparsification with Global Sparsity Constraint

Below, we present our proposed network sparsification framework and the method for solving the minimization problem in the framework.

3.1. A Probabilistic Sparsification Framework

Let \mathcal{D} be a dataset consisting of N i.i.d. samples $\{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$, $\mathbf{w} \in \mathbb{R}^n$ be the weights of a neural network. We denote $\mathbf{m} \in \{0, 1\}^n$ to be the masks of the weights. $m_i = 0$ means the weight w_i is pruned and otherwise w_i is kept. The problem of training sparse neural networks can be naturally formulated into the following empirical risk minimization problem:

$$\begin{aligned} \min_{\mathbf{w}, \mathbf{m}} \mathcal{L}(\mathbf{w}, \mathbf{m}) &:= \frac{1}{N} \sum_{i=1}^N \ell(h(\mathbf{x}_i; \mathbf{w} \circ \mathbf{m}), \mathbf{y}_i) \quad (1) \\ \text{s.t. } \mathbf{w} &\in \mathbb{R}^n, \|\mathbf{m}\|_1 \leq K \text{ and } \mathbf{m} \in \{0, 1\}^n, \end{aligned}$$

where $h(\cdot; \mathbf{w} \circ \mathbf{m})$ is output of the pruned network with \circ being the element-wise product of two vectors, and $\ell(\cdot, \cdot)$ is the loss function, e.g, the squared loss for regression or cross entropy loss for classification. $K = kn$ is the model size we want to reduce the network to, i.e., the number of remaining weights after pruning and k is the remaining ratio of model weights. In this framework, the model size is controlled by a single constraint which avoids tuning the pruning rate for each layer. However, since the objective is discrete with respect to the mask \mathbf{m} , problem (1) is hard to solve and thus cannot be applied in practice.

We notice that if we view each component of mask \mathbf{m} as a binary random variable and reparameterize problem (1) with respect to the distributions of this random variable, then problem (1) can be relaxed into an expected loss minimization problem over the weight and probability spaces, which is continuous. We need to point out that this is a very tight relaxation since empirical observations show that probabilities s_i converge to 0 or 1 after training (Section 5).

Specifically, we can view m_i as a Bernoulli random variable with probability s_i to be 1 and $1 - s_i$ to be 0, that is $m_i \sim \text{Bern}(s_i)$, where $s_i \in [0, 1]$. Assuming the variables m_i are independent, then we can get the distribution function of \mathbf{m} , i.e., $p(\mathbf{m}|\mathbf{s}) = \prod_{i=1}^n (s_i)^{m_i} (1-s_i)^{(1-m_i)}$. Thus, the model size can be controlled by the sum of the probabilities s_i , i.e., $\mathbf{1}^\top \mathbf{s}$, since $\mathbb{E}_{\mathbf{m} \sim p(\mathbf{m}|\mathbf{s})} \|\mathbf{m}\|_0 = \sum_{i=1}^n s_i$. Then the discrete constraint $\|\mathbf{m}\|_1 \leq K$ in problem (1) can be transformed into $\mathbf{1}^\top \mathbf{s} \leq K$ with each $s_i \in [0, 1]$, which is continuous and convex. Therefore, problem (1) can be relaxed into the following expected loss minimization problem:

$$\begin{aligned} \min_{\mathbf{w}, \mathbf{s}} \mathbb{E}_{p(\mathbf{m}|\mathbf{s})} \mathcal{L}(\mathbf{w}, \mathbf{m}) \quad (2) \\ \text{s.t. } \mathbf{w} \in \mathbb{R}^n, \mathbf{1}^\top \mathbf{s} \leq K \text{ and } \mathbf{s} \in [0, 1]^n. \end{aligned}$$

Discussion. Appealing features of ProbMask:

- The constraints in problem (2) can be rewritten as $\|\mathbf{s}\|_1 \leq K$ and $\mathbf{s} \in [0, 1]^n$. Due to this ℓ_1 norm constraint, the optimal \mathbf{s} is sparse. Most of s_i would be either 0 or 1 with high probability, making \mathbf{m} converge to a deterministic mask. Therefore, \mathbf{s} after training would have a quite low variance and thus the loss of a randomly sampled \mathbf{m} would be close to the expected loss in Eqn.(2).
- Compared with problem (1), problem (2) is continuous. Moreover, the feasible region of problem (2) is quite simple, which is actually the intersection of the cube $[0, 1]^n$ and the half space $\mathbf{1}^\top \mathbf{s} \leq K$. For such simple set, the projection operator has an explicit expression, please see Theorem 1 for the details. This makes it possible to adopt the efficient optimization algorithms such as projected gradient descent to solve problem (2).
- In our framework, the problem is reparameterized with respect to probability, which can be used as a global criterion to measure the importance of the weights in different layers. Note that the constraint is applied over all probability for different layers, rather than setting a uniform sparsity across layers. The amount of redundancy in each layer of the neural network can be learned automatically in the process of solving problem (2), which will be verified in Section 5. Therefore, we avoid setting pruning ratio for each layer manually. The benefits of globally comparable property of probability on the model size and accuracy will be further verified in Section 4.

3.2. Optimization with Projected Gradient Descent

Below, we present our training method for problem (2). We update both the weights \mathbf{w} and the probability \mathbf{s} at training time. At test time, we obtain the sparse network $\mathbf{w} \circ \mathbf{s}$

by sampling according to probability \mathbf{s} . We adopt projected gradient descent (PGD) as the optimizer and the details are as follows.

[Gradient Computation] The difficulty lies in computing the gradient of the expected loss with respect to the probability. Therefore, in this paper, we adopt Gumbel-Softmax [16, 25] trick to calculate the gradient, with which the gradient w.r.t. weights and probability can be calculated in the following form:

$$\begin{aligned} \nabla_{\mathbf{s}, \mathbf{w}} \mathbb{E}_{p(\mathbf{m}|\mathbf{s})} \mathcal{L}(\mathbf{w}, \mathbf{m}) \\ = \mathbb{E}_{\mathbf{g}_0, \mathbf{g}_1} \nabla_{\mathbf{s}, \mathbf{w}} \mathcal{L}\left(\mathbf{w}, \mathbf{1}\left(\log\left(\frac{\mathbf{s}}{\mathbf{1}-\mathbf{s}}\right) + \mathbf{g}_1 - \mathbf{g}_0 \geq 0\right)\right) \quad (3) \end{aligned}$$

$$\approx \mathbb{E}_{\mathbf{g}_0, \mathbf{g}_1} \nabla_{\mathbf{s}, \mathbf{w}} \mathcal{L}\left(\mathbf{w}, \sigma\left(\frac{\log\left(\frac{\mathbf{s}}{\mathbf{1}-\mathbf{s}}\right) + \mathbf{g}_1 - \mathbf{g}_0}{\tau}\right)\right) \quad (4)$$

$$\approx \frac{1}{I} \sum_{i=1}^I \nabla_{\mathbf{s}, \mathbf{w}} \mathcal{L}\left(\mathbf{w}, \sigma\left(\frac{\log\left(\frac{\mathbf{s}}{\mathbf{1}-\mathbf{s}}\right) + \mathbf{g}_1^{(i)} - \mathbf{g}_0^{(i)}}{\tau}\right)\right), \quad (5)$$

where $\mathbf{1}(\mathbf{A}) \in \{0, 1\}^n$ is the indicator function. \mathbf{g}_0 and \mathbf{g}_1 are two random variables in \mathbb{R}^n , with each element i.i.d sampled from Gumbel(0, 1) distribution. $\mathbf{g}_1^{(i)}$ and $\mathbf{g}_0^{(i)}$ with $i = 1, 2, \dots, I$ are $2I$ sampled instances. $\sigma(\cdot) : \mathbb{R}^n \rightarrow (0, 1)^n$ here is the element-wise sigmoid function, i.e., $\sigma(\mathbf{x}) = \frac{1}{1+e^{-\mathbf{x}}}$ for any $\mathbf{x} \in \mathbb{R}^n$. τ is a temperature annealing parameter decreasing linearly during training and precise choice of the decreasing function contributes to convergence of probability to a deterministic state. We will present empirical observations in Section 5 and provide some informal insights on such contribution from precise choice of temperature decreasing function in appendix. From Eqn.(4) to Eqn.(5), multiple networks are sampled to obtain a steady gradient flow with a low variance. The proof of the equations above is placed in appendix.

[Projected Gradient Descent] We denote the feasible region of probability in problem (2) as C , that is $C = \{\mathbf{s} \mid \|\mathbf{s}\|_1 \leq K \text{ and } \mathbf{s} \in [0, 1]^n\}$. The theorem below shows that the projection of a vector onto C can be calculated efficiently, which makes PGD applicable.

Theorem 1. For each vector \mathbf{z} , its projection \mathbf{s} in the set C can be calculated as follows:

$$\mathbf{s} = \min(1, \max(0, \mathbf{z} - v_2^* \mathbf{1})). \quad (6)$$

where $v_2^* = \max(0, v_1^*)$ with v_1^* being the solution of the following equation

$$\mathbf{1}^\top [\min(1, \max(0, \mathbf{z} - v_1^* \mathbf{1}))] - K = 0. \quad (7)$$

The equation (7) can be solved by bisection method efficiently. Now we can apply PGD to solve problem (2) directly on probability space with explicit sparsity constraint. We provide a complete view of ProbMask in Algorithm 1, and supplementary messages can be found in appendix.

Algorithm 1 Probabilistic Masking (ProbMask)

Input: target remaining ratio k_f , a dense network w .

- 1: Initialize w , assign probabilities s to weights w , let $s = \mathbf{1}$ and $\tau = k = 1$.
 - 2: **for** training epoch $t = 1, 2, \dots, T$ **do**
 - 3: Decrease the temperature annealing parameter by $\tau = 0.97(1 - t/T) + 0.03$.
 - 4: Update k according to Eqn.(8).
 - 5: **for** each training iteration **do**
 - 6: Sample mini batch of data $\mathcal{B} = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_B, \mathbf{y}_B)\}$.
 - 7: Generate $\mathbf{g}_1^{(i)}$ and $\mathbf{g}_0^{(i)}$ with each element sampled from Gumbel(0, 1), $i = 1, 2, \dots, I$.
 - 8: $\mathbf{s} \leftarrow \text{proj}_C(\mathbf{z})$, with $\mathbf{z} = \mathbf{s} - \eta \frac{1}{I} \sum_{i=1}^I \nabla_{\mathbf{s}} \mathcal{L}_{\mathcal{B}} \left(\mathbf{w}, \sigma \left(\frac{\log(\frac{\mathbf{s}}{1-\mathbf{s}}) + \mathbf{g}_1^{(i)} - \mathbf{g}_0^{(i)}}{\tau} \right) \right)$.
 - 9: $\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{1}{I} \sum_{i=1}^I \nabla_{\mathbf{w}} \mathcal{L}_{\mathcal{B}} \left(\mathbf{w}, \sigma \left(\frac{\log(\frac{\mathbf{s}}{1-\mathbf{s}}) + \mathbf{g}_1^{(i)} - \mathbf{g}_0^{(i)}}{\tau} \right) \right)$
 - 10: **end for**
 - 11: **end for**
 - 12: **return** A pruned network $w \circ m$ by sampling a mask m from the distribution $p(m|s)$.
-

[Gradually Increasing Pruning Rate] We increase the pruning rate gradually to make a smooth transformation from dense to sparse status. We utilize the increasing function of [41],

$$k = k_f + (1 - k_f) \left(1 - \frac{t - t_1}{t_2 - t_1} \right)^3, \quad (8)$$

where $t \in \{t_1, t_1 + 1, \dots, t_2\}$ is the current epoch number and k_f is the targeted remaining ratio. k keeps 1 before epoch t_1 and k_f after epoch t_2 .

Remark 1. *ProbMask directly works on the probability space without any further reparameterization, avoiding the drawback of gradient vanishing [32, 27]. Together with the global sparsity constraint, ProbMask finally learns a deterministic state of probability, resulting in a little performance gap in testing and training phases.*

Remark 2. *ProbMask can be trained with randomly initialized weights or from pretrained weights. ProbMask can explicitly control the sparsity by choosing a proper K to achieve a desired model size and does not need to search any parameters.*

4. Experiment

In this section, we conduct a series of experiments to evaluate the performance of our proposed method. We divide the experiments into two parts. In part one, we conduct lots of relatively small-scaled experiments on CIFAR-10/100 datasets with modern architectures VGG19 [31] and ResNet32 [13] to verify some appealing properties of our method. In part two, we verify the superiority of our method over state-of-the-art methods by conducting experiments on ImageNet [3]. We choose six representative methods PBW (Pruning by Weight, [10]), MLPrune [39], RIGL [5], STR [18], DNW[25], GMP [41]) as baselines. PBW [10] is a

classic magnitude-based pruning method. MLPrune [39] is a latest Hessian-based pruning method showing overall better performance [33] against various sparse-to-sparse training methods (SET [26], DEEPR [1], DSR [28]), so we compare with these sparse-to-sparse training methods implicitly in CIFAR experiments. DNW [35], GMP [41], STR [18] are state-of-the-art methods on dense-to-sparse training. RIGL [5] is the state-of-the-art sparse-to-sparse training method. Due to the space limitation, we postpone the experimental configurations and MobileNetV1 [15] experiments into appendix.

4.1. VGG19 and ResNet32 on CIFAR-10/100

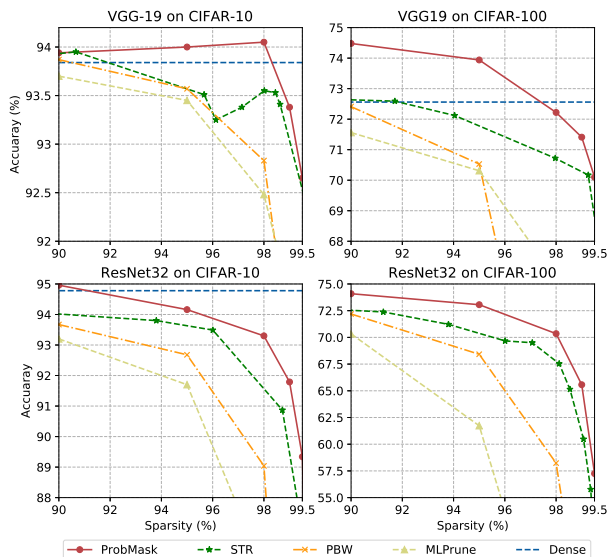


Figure 1. Comparison of Top-1 Accuracy on CIFAR-10/100.

Table 1 presents the detailed accuracy of PBW, MLPrune

and ProbMask at different pruning ratios. It is very hard to accurately tune weight decay parameter in STR to obtain the desired pruning ratio. Therefore we tune the weight decay parameter manually to make it have roughly the same pruning ratio range with ProbMask, i.e., 90% to 99.9%.

The results in both Table 1 and Figure 1 demonstrate that our ProbMask can steadily outperform the baselines and the superiority becomes more significant at higher pruning ratios. From Table 1, we can see that when prune rate come to 99.5% or higher on CIFAR-10/100, PBW and ML-Prune would seriously degrade or even collapse, while our ProbMask can still achieve significantly higher. Figure 1 shows that on CIFAR-100 with VGG19, the gap between ProbMask and STR would be roughly 2% on average when the remaining ratio is in the range of [0.9, 0.98]. Pruning ResNet32 is more challenging since VGG19 has about 10 times parameters than ResNet32. In this case, the gap becomes more significant especially at high pruning ratios, which can be up to 5% on CIFAR-100 experiments. The superiority of ProbMask over such high prune ratios attributes to our global sparsity constraint, allowing us to have non-uniform sparsity budgets across layers. This will be further validated in the ablation study in Section 5.

4.2. ResNet50 on ImageNet-1K

In this section, we evaluate the performance of our ProbMask on ImageNet with ResNet50. Table 2 and Figure 2 report the detailed accuracy at different pruning ratios. ProbMask steadily outperforms state-of-the-art methods with a large margin, especially when the pruning ratio is high than 98%. Notably the gap comes up to 5% at 98% sparsity and 10% at 99% sparsity. DNW and GMP allocate uniform sparsity budget. They present reasonably good performance at 90% sparsity while falling behind ProbMask by about 9% at 98% sparsity. This validates our previous claim that identifying weight allocation for different layers really matters. Uniform sparsity budget is a reasonable compromise but obviously don't give a perfect solution. STR attempts to learn weight allocation for different layers but don't give perfect results. ProbMask presents much better performance on high sparsity regions, leading to a gap about 10% percent at 99% sparsity. With the global comparable nature of probability, ProbMask easily learns a much better weight allocation scheme for different layer. We also compare ProbMask with Sparse VD [27] on sparsity 90%. Sparse VD finds a subnet with 73.84% Top-1 Accuracy, a weaker result than ProbMask. We also observe noticeable fluctuations between different runs, and this can be expected because Sparse VD adopts crude cut-off practice rather than sampling. This inevitably results in performance gap in training and testing phases. ProbMask learns a deterministic mask at the end of training, fixing the training and testing performance discrepancy problem. Figure 3 reports the accuracy-versus-

FLOPs for ProbMask and compared methods. It shows that ProbMask finds a smaller mask with comparable accuracy and FLOPs and achieve state-of-the-art result on it.

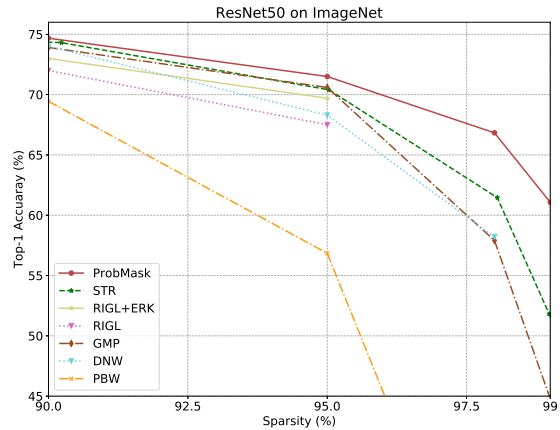


Figure 2. ProbMask comfortably beats state-of-the-art methods in all sparsity regions. Notably, the gap comes up to 5% at 98% sparsity and 10% at 99% sparsity.

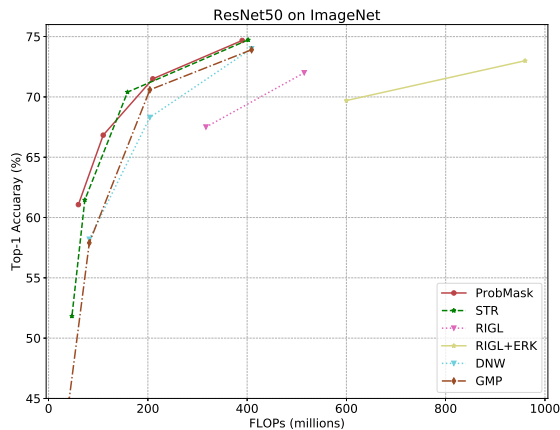


Figure 3. ProbMask obtains a smaller sparse network with comparable accuracy and FLOPs, still achieving state-of-the-art result on accuracy-versus-FLOPs curve.

4.3. Powerful Tool for Finding Supermasks

Previous works on supermasks, i.e, subnetworks achieving good performance with weights fixed at random state, focus on sparsity region [10%, 90%]. Here, we would like to explore the performance of supermasks with higher sparsity, [90%, 99%]. We conduct experiments on modern architecture ResNet32 and dataset CIFAR-100, a harder task than CIFAR-10 where a large portion of previous experiments are conducted. In this experiment, weights are fixed at initialization state by Kaiming Normal [12]. Hyperparameters follow the same as previous CIFAR experiments. According to Figure 4, we observe that ProbMask easily scales to ultra sparse region with about 50% accuracy and

Dataset	CIFAR-10						CIFAR-100					
	90%	95%	98%	99%	99.5%	99.9%	90%	95%	98%	99%	99.5%	99.9%
VGG19	93.84	-	-	-	-	-	72.56	-	-	-	-	-
PBW [9]	93.87	93.57	92.83	90.89	10.00	10.00	72.41	70.53	58.91	1.00	1.00	1.00
MLPrune [39]	93.70	93.45	92.48	91.44	88.18	65.38	71.56	70.31	66.77	60.10	50.98	5.58
ProbMask	93.94	94.00	94.05	93.38	92.65	89.79	74.48	73.94	72.22	71.41	70.10	60.41
ResNet32	94.78	-	-	-	-	-	75.94	-	-	-	-	-
PBW [9]	93.67	92.68	89.04	77.03	73.03	38.64	72.19	68.42	58.23	43.00	20.75	5.96
MLPrune [39]	93.20	91.70	85.64	76.88	67.66	36.09	70.33	61.73	37.86	22.38	13.85	5.50
ProbMask	94.96	94.16	93.30	91.79	89.34	76.87	74.09	73.06	70.35	65.57	57.25	26.72

Table 1. Accuracy of VGG19 and ResNet32 on CIFAR-10/100 at different pruning ratios.

2% remaining weights, while state-of-the-art method edge-popup [29] collapse with less than 30% accuracy. It is a surprising result that a subnet with 2% fixed random weights still succeeds in obtaining nearly 50% accuracy on a task with 100 categories. It shows that the structure in networks already provides valuable information for classification.

Dataset	ImageNet			
	90%	95%	98%	99%
ResNet50	77.01	-	-	-
PBW[9]	69.44	56.84	22.46	5.98
MLPrune[39]	60.98	30.89	3.16	0.77
GMP[41]	73.91	70.59	57.90	44.78
DNW[35]	74.00	68.30	58.20	-
STR[18]	74.31	70.40	61.46	50.35
RIGL[5]	72.00	67.50	-	-
ProbMask	74.68	71.50	66.83	61.07

Table 2. Accuracy of ResNet50 on ImageNet at different pruning ratios. ProbMask steadily beats previous state-of-the-art methods on Hessian-based pruning, weight magnitude pruning, dense-to-sparse training and sparse-to-sparse training. RIGL improves with the help of ERK (Erdős-Rényi-Kernel) but will result in doubling the FLOPs at inference time, so we put it in Figure 2).

5. Further Analysis

[Global Comparability of Probability] Table 1 shows that PBW and MLPrune collapse on CIFAR-10/100 when the pruning ratio is as high as 99.9%. To explore the reason, we plot the remaining ratio across layers at pruning ratio of 90% and 99.9% on CIFAR-10 in Figure 5. It shows that when the pruning ratio is high, PBW and MLPrune prune almost all the weights in certain layers with remaining ratio approaching 10^{-6} . The reason is that the proposed weight importance measure in PBW and MLPrune are not glob-

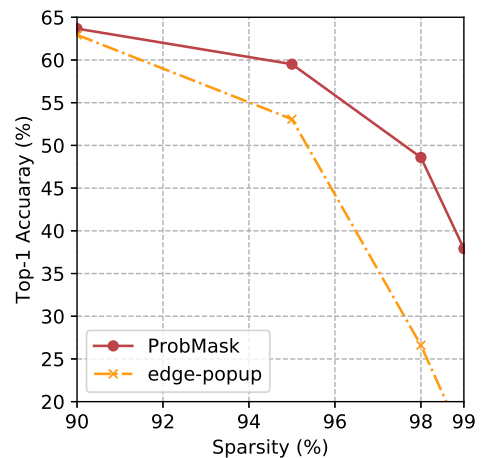


Figure 4. ProbMask can find a supermask with just 2% remaining weights and nearly 50% accuracy on CIFAR-100. Weights are fixed at initialization state.

ally comparable. Although the weight importance scores in different layers have been normalized in MLPrune, their magnitudes are still quite different. A global threshold could remove almost all the weights in certain layers in order to achieve high enough pruning ratio. The pruning ratio of ProbMask varies in a proper range, attributed to the global comparable nature of probability. We observe that the learned sparsity budget is a wise balance between uniform sparsity budget and cutting one layer off. The first and last layer are assigned a bit more budget above average and several important bottleneck layers are detected automatically to assign more budget.

[Superiority of Global Sparsity Constraint over Layer-wise Constraint] In layer-wise constraint, we force all the pruning ratios in each layer to be equal and also equal to the one in the global constraint. The experiment is conducted on CIFAR-10 with ResNet32 and the results are given in Table 3. It shows that the gap grows up rapidly when the pruning ratio is larger than 98%. For example,

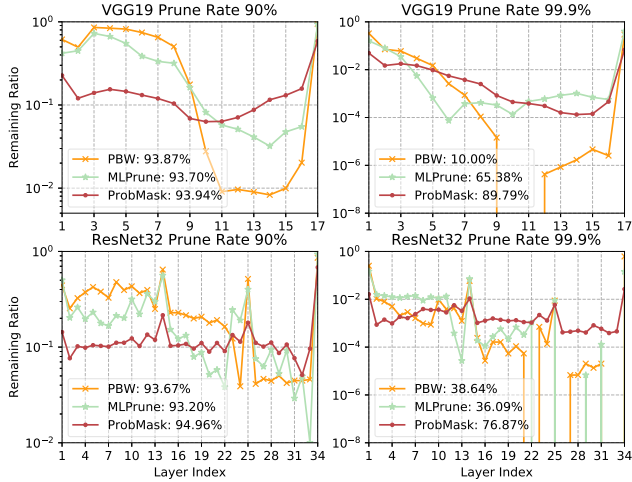


Figure 5. ProbMask learns a wise balance between uniform sparsity budget and abominably cutting one layer off, leading to compelling performance over sparsity range [90%, 99.9%].

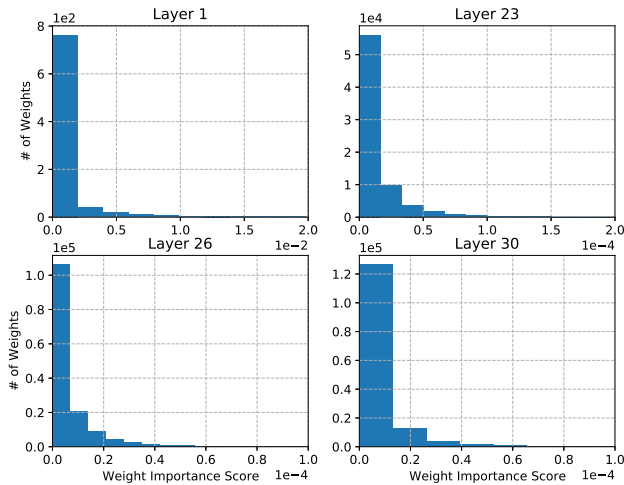


Figure 6. Weight importance score histogram of ResNet32 from MLPrune with pruning rate 99.9%. Note that the index of x-axis is scaled to $1e-2$ for Layer 1, and $1e-4$ for Layer 23, 26, 30. This means that there exist two orders of magnitude difference across layers among weight importance scores.

when the pruning ratio is 99.9%, the accuracy of global sparsity constraint can be up to 57.75% higher than the layer-wise one. This points out the importance of identifying sparsity budget for different layers again by ablation study.

[Convergence to Deterministic Mask] To show that the mask trained by our ProbMask can converge to a deterministic mask after training, we randomly choose some layers from VGG19 and present their distribution of the probability value after training in Figure 7. We can see that after training, almost all of the probabilities s_i can converge to either 0 and 1, leading to a deterministic mask. This attributes

Dataset	CIFAR-10					
Ratio	90%	95%	98%	99%	99.5%	99.9%
ResNet32	94.78	-	-	-	-	-
LSB	94.89	94.09	92.64	90.89	74.6	19.12
GSB	94.96	94.16	93.30	91.79	89.34	76.87

Table 3. Comparing Layerwise Sparsity Budget (LSB, assigning uniform budget across layers) and Global Sparsity Budget (GSB) of ProbMask on ResNet32. GBS begins to take effect when sparsity comes up to 98% and becomes prominent when sparsity is larger than 99.5%.

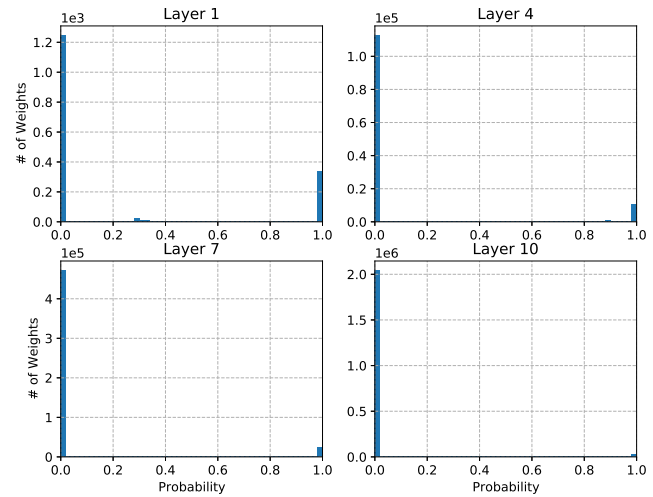


Figure 7. Probability histogram of VGG19 trained by ProbMask on CIFAR-10 at pruning rate 90%.

to ℓ_1 norm in our global sparsity constraint over the probability space and the precise chosen temperature annealing scheme.

6. Conclusion

This paper proposes an effective network sparsification method ProbMask and demonstrates state-of-the-art results on various models and datasets. We provide evidence that probability can serve as a suitable global comparator to measure weight importance and solve the training and testing performance discrepancy problem observed in practice. ProbMask can also serve as a powerful tool for identifying subnetworks with high performance in a randomly weighted dense neural network.

Acknowledgements

This work is supported by GRF 16201320.

References

- [1] Guillaume Bellec, David Kappel, Wolfgang Maass, and Robert Legenstein. Deep rewiring: Training very sparse deep networks. *arXiv preprint arXiv:1711.05136*, 2017. [2](#), [3](#), [5](#)
- [2] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013. [3](#)
- [3] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. [5](#)
- [4] Tim Dettmers and Luke Zettlemoyer. Sparse networks from scratch: Faster training without losing performance. *arXiv preprint arXiv:1907.04840*, 2019. [2](#), [3](#)
- [5] Utku Evci, Trevor Gale, Jacob Menick, Pablo Samuel Castro, and Erich Elsen. Rigging the lottery: Making all tickets winners. In *International Conference on Machine Learning*, pages 2943–2952. PMLR, 2020. [2](#), [3](#), [5](#), [7](#)
- [6] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018. [2](#), [3](#), [11](#)
- [7] Trevor Gale, Erich Elsen, and Sara Hooker. The state of sparsity in deep neural networks. *arXiv preprint arXiv:1902.09574*, 2019. [1](#), [2](#), [3](#)
- [8] Yiwen Guo, Anbang Yao, and Yurong Chen. Dynamic network surgery for efficient dnns. In *Advances in neural information processing systems*, pages 1379–1387, 2016. [1](#), [2](#)
- [9] Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015. [1](#), [2](#), [7](#)
- [10] Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. In *Advances in neural information processing systems*, pages 1135–1143, 2015. [2](#), [5](#)
- [11] Babak Hassibi and David G Stork. Second order derivatives for network pruning: Optimal brain surgeon. In *Advances in neural information processing systems*, pages 164–171, 1993. [2](#)
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015. [6](#)
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [5](#)
- [14] Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1389–1397, 2017. [2](#)
- [15] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. [5](#), [11](#)
- [16] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016. [4](#)
- [17] Minsoo Kang and Bohyung Han. Operation-aware soft channel pruning using differentiable masks. *arXiv preprint arXiv:2007.03938*, 2020. [2](#)
- [18] Aditya Kusupati, Vivek Ramanujan, Raghav Somani, Mitchell Wortsman, Prateek Jain, Sham Kakade, and Ali Farhadi. Soft threshold weight reparameterization for learnable sparsity. *arXiv preprint arXiv:2002.03231*, 2020. [2](#), [3](#), [5](#), [7](#), [11](#)
- [19] Yann LeCun, John S Denker, and Sara A Solla. Optimal brain damage. In *Advances in neural information processing systems*, pages 598–605, 1990. [2](#)
- [20] Namhoon Lee, Thalaiyasingam Ajanthan, and Philip HS Torr. Snip: Single-shot network pruning based on connection sensitivity. *arXiv preprint arXiv:1810.02340*, 2018. [2](#), [3](#)
- [21] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. *arXiv preprint arXiv:1608.08710*, 2016. [1](#), [2](#)
- [22] Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. Rethinking the value of network pruning. *arXiv preprint arXiv:1810.05270*, 2018. [1](#), [2](#), [11](#)
- [23] Christos Louizos, Max Welling, and Diederik P Kingma. Learning sparse neural networks through l_0 regularization. *arXiv preprint arXiv:1712.01312*, 2017. [2](#), [3](#)
- [24] Jian-Hao Luo, Jianxin Wu, and Weiyao Lin. Thinet: A filter level pruning method for deep neural network compression. In *Proceedings of the IEEE international conference on computer vision*, pages 5058–5066, 2017. [2](#)
- [25] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016. [4](#), [5](#)
- [26] 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):1–12, 2018. [2](#), [3](#), [5](#)
- [27] Dmitry Molchanov, Arsenii Ashukha, and Dmitry Vetrov. Variational dropout sparsifies deep neural networks. *arXiv preprint arXiv:1701.05369*, 2017. [2](#), [3](#), [5](#), [6](#)
- [28] Hesham Mostafa and Xin Wang. Parameter efficient training of deep convolutional neural networks by dynamic sparse reparameterization. *arXiv preprint arXiv:1902.05967*, 2019. [2](#), [3](#), [5](#)
- [29] Vivek Ramanujan, Mitchell Wortsman, Aniruddha Kembhavi, Ali Farhadi, and Mohammad Rastegari. What’s hidden in a randomly weighted neural network? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11893–11902, 2020. [3](#), [7](#), [11](#)
- [30] Alex Renda, Jonathan Frankle, and Michael Carbin. Comparing rewinding and fine-tuning in neural network pruning. *arXiv preprint arXiv:2003.02389*, 2020. [2](#)

- [31] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. [5](#)
- [32] Suraj Srinivas, Akshayvarun Subramanya, and R Venkatesh Babu. Training sparse neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 138–145, 2017. [2](#), [3](#), [5](#)
- [33] Chaoqi Wang, Guodong Zhang, and Roger Grosse. Picking winning tickets before training by preserving gradient flow. *arXiv preprint arXiv:2002.07376*, 2020. [3](#), [5](#), [11](#)
- [34] Ziheng Wang, Jeremy Wohlwend, and Tao Lei. Structured pruning of large language models. *arXiv preprint arXiv:1910.04732*, 2019. [2](#)
- [35] Mitchell Wortsman, Ali Farhadi, and Mohammad Rastegari. Discovering neural wirings. In *Advances in Neural Information Processing Systems*, pages 2684–2694, 2019. [2](#), [5](#), [7](#)
- [36] Xia Xiao, Zigeng Wang, and Sanguthevar Rajasekaran. Autoprune: Automatic network pruning by regularizing auxiliary parameters. In *Advances in Neural Information Processing Systems*, pages 13681–13691, 2019. [2](#), [3](#)
- [37] Sirui Xie, Hehui Zheng, Chunxiao Liu, and Liang Lin. Snas: stochastic neural architecture search. *arXiv preprint arXiv:1812.09926*, 2018. [11](#)
- [38] Mao Ye, Chengyue Gong, Lizhen Nie, Denny Zhou, Adam Klivans, and Qiang Liu. Good subnetworks provably exist: Pruning via greedy forward selection. *arXiv preprint arXiv:2003.01794*, 2020. [2](#)
- [39] Wenyuan Zeng and Raquel Urtasun. MLPrune: Multi-layer pruning for automated neural network compression, 2019. [1](#), [2](#), [5](#), [7](#)
- [40] Hattie Zhou, Janice Lan, Rosanne Liu, and Jason Yosinski. Deconstructing lottery tickets: Zeros, signs, and the supermask. In *Advances in Neural Information Processing Systems*, pages 3597–3607, 2019. [2](#), [3](#)
- [41] Michael Zhu and Suyog Gupta. To prune, or not to prune: exploring the efficacy of pruning for model compression. *arXiv preprint arXiv:1710.01878*, 2017. [2](#), [5](#), [7](#), [11](#)