# Pruning-aware Sparse Regularization for Network Pruning

Nanfei Jiang[1,2,*], Xu Zhao[1,*], Chaoyang Zhao[1], Yongqi An[1,2], Ming Tang[1,2], Jinqiao Wang[1,2]

1. National Laboratory of Pattern Recognition, Institute of Automation,
Chinese Academy of Sciences,Beijing, 100190, China.
2. University of Chinese Academy of Sciences, Beijing, 100049, China

{nanfei.jiang, xu.zhao, chaoyang.zhao, yongqi.an, tangm, jqwang}@nlpr.ia.ac.cn

## Abstract

*Structural neural network pruning aims to remove the redundant channels in the deep convolutional neural networks (CNNs) by pruning the filters of less importance to the final output accuracy. To reduce the degradation of performance after pruning, many methods utilize the loss with sparse regularization to produce structured sparsity. In this paper, we analyze these sparsity-training-based methods and find that the regularization of unpruned channels is unnecessary. Moreover, it restricts the network's capacity, which leads to under-fitting. To solve this problem, we propose a novel pruning method, named MaskSparsity, with pruning-aware sparse regularization. MaskSparsity imposes the fine-grained sparse regularization on the specific filters selected by a pruning mask, rather than all the filters of the model. Before the fine-grained sparse regularization of MaskSparity, we can use many methods to get the pruning mask, such as running the global sparse regularization. MaskSparsity achieves 63.03%-FLOPs reduction on ResNet-110 by removing 60.34% of the parameters, with no top-1 accuracy loss on CIFAR-10. On ILSVRC-2012, MaskSparsity reduces more than 51.07% FLOPs on ResNet-50, with only a loss of 0.76% in the top-1 accuracy.*

*The code is released at https://github.com/ CASIA-IVA-Lab/MaskSparsity. Moreover, we have integrated the code of MaskSparity into a PyTorch pruning toolkit, EasyPruner, at https://gitee.com/ casia_iva_engineer/easypruner.*

## 1. Introduction

Convolutional Neural Networks (CNNs) have demonstrated a great success on a variety of computer vision tasks, like image classification [30], detection [25], and semantic segmentation [4]. However, the increasing depth and width of
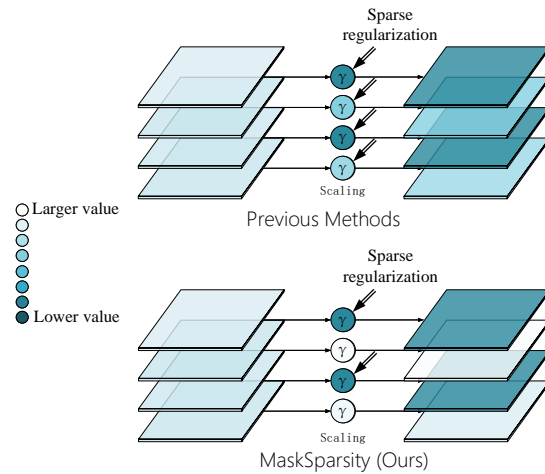
---

*Equal contribution.



Figure 1. Visual comparison between the previous sparsity-training-based methods and the proposed MaskSparsity method. MaskSparity apply the sparse regularization only on the scaling factors of less-important channels.

the CNNs also lead to higher computing resources demands and excessive memory footprint requirements. Typically, the widely used ResNet models [8] have millions of parameters, requiring billions of float point operations (FLOPs), making it a great challenge to deploy most state-of-the-art CNNs on edge devices. Network pruning is an effective way to compress and accelerate CNNs. It is attracting much attention from researchers. It can remove the parameters in the deep CNNs and reduce the required FLOPs and memory footprint while preserving the performance.

A typical scheme of network pruning consists of three stages: (1) training an over-parameterized model normally; (2) pruning the model under a certain criterion; and (3) fine-tuning the pruned model to reduce the degradation caused by pruning. Some of the existing network pruning methods apply a sparsity training stage after step (1). These methods apply sparse regularization on the filter weights of the convolution layers [1, 34] or scaling factors [14, 26] of the
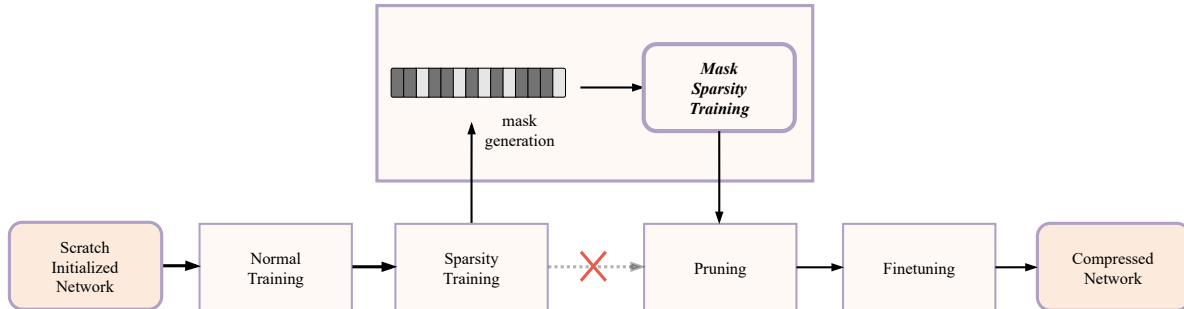
Figure 2. The pipeline of the proposed MaskSparsity method.

batch normalization layers. After the sparsity training, the corresponding filter weights or scaling factors of unimportant channels are considered to be near zero. Then these channels could be safely pruned without affecting the output values of the corresponding layers too much. We call these methods sparsity-training-based methods.

In the sparsity-training-based methods, to get an expected sparse rate of the model, they adopt the global sparse regularization. However, in these methods, the weights of important channels are regularized in the sparsity training stage, although they are preserved after the pruning. It is generally regarded that proper regularization achieved a good result by avoiding over-fitting. However, the model will be under-fitted when the regularization coefficient is too large. Since the weights of important channels are also regularized by the sparse regularization, the magnitude of these weights is usually decayed towards 0. This prevents the coverage to better local minima of the network in the sparsity training stage, which affects the final performance of the fine-tuned pruned network. Moreover, the sparse rate of the globally trained network is hard to control. This usually results in the inconsistency between the sparse mask and the pruning mask if we want to prune a model to a pre-defined FLOPS.

To address the problem mentioned above, we propose a novel sparsity-training-based channel pruning approach, MaskSparsity. Different from the previous sparsity-training-based methods which impose the regularization on all channels of each layer, MaskSparsity only imposes the regularization on the specific channels selected by the pruning mask, which indicates the unimportant channels, as shown in Figure 1. Through this mask, MaskSparsity can realize the strong correlation between pruning and regularization, and carry out a pruning-aware regularization. In other words, we only impose regularization on the channels to be pruned and prune the channels where regularization is applied. The perfect match between the sparse channels and the pruning channels allows us to minimize the impact of sparse regularization and maximize the accuracy of the pruned networks.

Compared with the typical pruning methods which di-

rectly prune the unimportant channels, the MaskSparsity can gradually push the unimportant parameters towards zero in a long period of iterations during the sparsity training stage. This prevents the model from a dramatic change of structure or weight in the pruning stage. It is regarded that the dramatic change may result in a certain amount of information loss which is harmful to restoring the accuracy of the model in the fine-tuning stage.

To summarize, our main contributions are three-fold:

- We analyze the previous sparsity-training-based methods in the previous work, which simply impose L1 regularization on all channels of the model. We find out the over-regularization problem on important channels.

- We propose MaskSparsity to solve these problems by more fine-grained sparsity training.

- The extensive experiments on two benchmarks show the effectiveness and efficiency of MaskSparsity.

## 2. Related Work

We mainly focus on the structural pruning methods in this paper. In this section, we first review the closely related works, *i.e.* the sparsity-training-based structural pruning methods. After that, we list other structural pruning methods.

### 2.1. Sparsity-training-based Pruning Methods

To make the network adaptively converge to a sparse structure and alleviate the damage of the pruning process to the network's output, some sparsity-training-based pruning methods are proposed. There are mainly two categories of these methods according to the place where sparse regularization is applied.

The first category of methods is the group-sparsity-based methods that apply spare regularization on the filter weights. Alvarez and Salzmann [1] proposed to use a group sparsity regularizer to determine the number of channels of each layer. Wen *et al.* [34] proposed a Structured

Sparsity Learning (SSL) method to regularize the structure to obtain a hardware-friendly pruned structure. Alvarez and Salzmann [2] added a low-rank regularizer to improve the pruning performance. Li and Gu proposed the Hinge [20] by combining the filter pruning and low-rank decomposition into the group sparsity training framework.

The second category of methods is the indirect group-sparse methods, which apply the sparse regularization on the scaling factors of each layer. The representative method is NetSlim [26] method, which sparsely regularizes the scaling factors of BN layers to get the sparse structure and remove less important channels. Huang and Wang [14] proposed to add a new scaling factor vector to each layer to apply the sparse regularization. Srinivas and Subramanya [32] proposed to impose sparse constraint over each weight with additional gate variables and achieve high compression rates by pruning connections with zero gate values. Ye and You [37] proposed to prune channels with layer-dependent thresholds according to the different weight distribution of each layer. [40] develop the norm-based importance estimation by taking the dependency between the adjacent layers into consideration.

These methods apply the global sparse regularization on the network channels, which over-regularize the important channels. Our MaskSparsity method solves this problem and can improve the performance of sparsity-training-based methods to the new state-of-the-art.

## 2.2. Non-sparsity-training-based Pruning Methods

Recently, many non-sparsity-training-based pruning methods also show good performance. These methods usually evaluate the importance of each channel with a hand-craft criterion first. After that, they directly prune the unimportant channels and finetune the network. For instance, Li and Kadav [18] proposed to prune filters with smaller L1 norm values in a network. Based on the theory of Geometric Median (GM) [5], He and Liu proposed FPGM [13] to prune the filters with the most replaceable contribution. Inspired by the discovery that the average rank of multiple feature maps generated by a single filter is always the same, Lin *et al.* [21] proposed to prune filters by exploring the High Rank of feature maps (HRank). In this paper, we compare the performance of the proposed MaskSparsity with these methods and show good pruning performance.

## 3. Methodology

### 3.1. Notations

We assume that a convolutional neural network consists of multiple convolutional layers and each convolution layer is followed by a batch normalization (BN) [15] layer. For the $l$-th convolutional layer, we use $C_l$ and $N_l$ to represent the number of its input channels and output channels, and

$k_l \times k_l$ represent the kernel size.

We use $\mathcal{W}^{(l)} = \{\mathcal{W}^{(l)}_{1,:,:,:}, \mathcal{W}^{(l)}_{2,:,:,:}, ..., \mathcal{W}^{(l)}_{N_l,:,:,:}\} \in \mathbb{R}^{N_l \times C_l \times k_l \times k_l}$ to represent the filters of the $l$-th convolutional layer. The input feature maps and the output feature maps to filters are denoted as $\mathcal{I}^{(l)} = \{\mathbf{i}^{(l)}_1, \mathbf{i}^{(l)}_2, ..., \mathbf{i}^{(l)}_{C_l}\} \in \mathbb{R}^{B \times C_l \times h_l \times w_l}$ and $\mathcal{O}^{(l)} = \{\mathbf{o}^{(l)}_1, \mathbf{o}^{(l)}_2, ..., \mathbf{o}^{(l)}_{N_l}\} \in \mathbb{R}^{B \times N_l \times h_l{}' \times w_l{}'}$. Here, $h_l$, $w_l$, $h_l{}'$ and $w_l{}'$ are the heights and widths of the input and output feature maps respectively. $B$ is the batch size of the input images. The $i$-th channel feature map $\mathbf{o}^{(l)}_i \in \mathbb{R}^{B \times h_l{}' \times w_l{}'}$ is generated by $\mathcal{W}^{(l)}_{i,:,:,:} \in \mathbb{R}^{C_l \times k_l \times k_l}$ and $\mathbf{I}^{(l)} \in \mathbb{R}^{B \times C_l \times h_l \times w_l}$.

For the $i$-th channel of the $l$-th BN layer with mean $\mu^{(l)}_i$, standard deviation $\sigma^{(l)}_i$, learned scaling factor $\gamma^{(l)}_i$ and bias $\beta^{(l)}_i$, regardless of bias of the convolutional layer, we have

$$\mathbf{o}^{(l)}_i = (\mathcal{W}^{(l)}_{i,:,:,:} \otimes \mathbf{I}^{(l)} - \mu^{(l)}_i)\frac{\gamma^{(l)}_i}{\sigma^{(l)}_i} + \beta^{(l)}_i. \tag{1}$$

### 3.2. Existing Sparsity-training-based Methods

Existing sparsity-training-based methods utilize sparse regularization loss to produce structured sparsity. Usually, the sparse regularization is either applied on the filter weights of convolutions or the channel scaling factors of BNs.

**(a) Sparsity on filter weights.** When sparse regularization is applied on the filter weights of convolutions, the training objective function of this category of methods is shown in Equation 2:

$$L_{\text{Sparsity}}(\mathcal{I}^0, y, \mathcal{W}) = L(f(\mathcal{I}^0, \mathcal{W}), y) + \lambda \cdot \sum_{l=1}^{L} \sum_{i=1}^{N_l} ||\mathcal{W}^{(l)}_{i,:,:,:}||_g, \tag{2}$$

where $(\mathcal{I}^0, y)$ denote the training samples and the labels, $\mathcal{W}$ denotes the trainable weights, the $L(f(\mathcal{I}^0, \mathcal{W}), y)$ is the objective function of normal training, $||\mathcal{W}^{(l)}_{i,:,:,:}||_g$ is a sparsity regularization penalty on the filter weights $\mathcal{W}$, here $|| \cdot ||_g$ is the group Lasso, $||\mathcal{W}^{(l)}_{i,:,:,:}||_g = \sqrt{\sum^{C_l} \sum^{k_1} \sum^{k_1} \left(\mathcal{W}^{(l)}_{i,j,k_1,k_2}\right)^2}$. $\lambda$ is the factor of controlling the strength of sparsity.

**(b) Sparsity on channel scaling factors.** When sparse regularization is applied on the channel scaling factors of the BN layer, the training objective function of this category of methods is shown in Equation 3:

$$L_{\text{Sparsity}}(\mathcal{I}^0, y, \mathcal{W}) = L(f(\mathcal{I}^0, \mathcal{W}), y) + \lambda \cdot \sum_{l=1}^{L} \sum_{i=1}^{N_l} ||\gamma^{(l)}_i||_g, \tag{3}$$

where $(\mathcal{I}^0, y)$, $L(f(\mathcal{I}^0, \mathcal{W}), y)$, and $\lambda$ denote the same mean as above. $||\gamma^{(l)}_i||_g$ is a sparsity regularization penalty

on the scaling factors $\gamma$ of BN layers. $||\cdot||_g$ is usually set as L1 regularization, and L2 regularization is available either as [7, 33].

In this paper, we choose the sparsity regularization penalty on **the channel scaling factors** for further investigation.

As with the normal regularization methods, the received gradients of the scaling factors from the normal training loss $L$ and the sparsity regularization $||\cdot||_g$ are usually against each other during training. The former aims at improving the model performance on the training set. The latter aims at restricting the range of the parameters and increasing the structure sparsity, which tends to increase the loss on the training set. The sparsity-training-based pruning method thinks that the unimportant convolutional channels are easily pushed to near 0 by $||\cdot||_g$, while the value of important channels is kept large by $||\cdot||_g$.

### 3.3. The Over-regularization Problem of Existing Sparsity-training-based Methods

Figure 3 shows the statistical results of two sets of scaling factors (absolute value) collected from the normally-trained and sparsely-trained ResNet-50 on ILSVRC-2012. In Figure 3, the purple histogram is the distribution of the normally-trained network, while the green histogram represents that of the sparse-trained network. Figure 3 shows that the scaling factors of the normally-trained network form one peak and those of the sparsely-trained network form two peaks. This is consistent with the bimodal-distribution observation of OT [38].
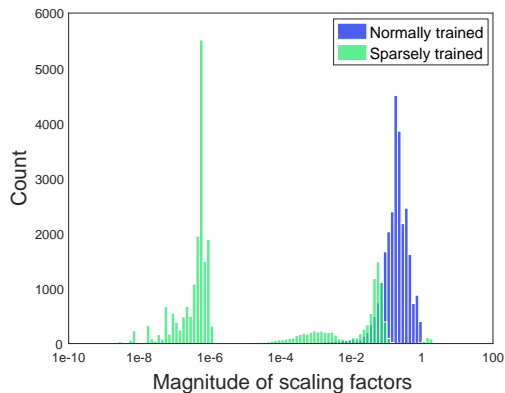


Figure 3. Distribution of scaling factors of ResNet50 before and after global sparsity training.

Obviously, the left scaling-factors peak of the sparsely-trained network represents the unimportant channels and the right peak represents the important channels. With Figure 3, we demonstrate the over-regularization problem of the existing sparsity-training-based methods, which are mentioned in the introduction.

It can be seen that the right peak of the sparsely-trained

network moves to 0 obviously, compared with their location in the histogram of the normally-trained network. This is a common phenomenon of the model regularization methods, i.e. the regularization cause a smaller magnitude of the model parameters. A small regularization usually leaders to better generalization performance on the test set. However, a too large regularization leads to under-fitting. This is because the regularization limits the network's capacity.

The modern CNNs are usually trained with *weight decay*, which is widely regarded to be similar to L2 regularization, especially under the SGD optimizer [6]. This is usually tuned to a proper magnitude to get the best performance. Moreover, the sparse regularization is regarded to push a large partition of a channel to be near 0. This requires a large weight for the sparsity loss. Therefore, we think the newly applied sparse regularization on the unpruned channels is over-regularization. It should be avoided.

### 3.4. Pruning-aware Sparse Regularization

Therefore, we propose a fine-grained sparsity training method that only applies the sparse regularization on the unimportant channels to keep a maximum representation ability of the important channels.

The task of sparsity training consists of two sub-tasks implicitly. The first sub-task is identifying the unimportant channels. The second sub-task is pushing the filter weights or scaling factors of unimportant channels to 0 by the sparse regularization loss. Existing sparsity-training-based methods accomplish the two sub-tasks simultaneously in the sparsity training stage. We propose to decouple the two sub-tasks. By doing this, we can apply the fine-grained spare regularization, which only sparse out the unimportant channels.

Figure 2 shows the training pipeline of the proposed MaskSparsity. We transform the sparsity training stage of the existing methods into two stages. The first stage is the sparsity training stage with global sparse regularization, which is aimed to get the indexes of the unimportant channels. The indexes are transformed into a binary pruning mask in previous methods. In our method, we use the mask to identify which channels to apply the sparse regularization in the second stage. To get the pruning mask, we directly threshold the scaling factors of the normally trained network. The details is shown in Equation 4:

$$\mathcal{M} = \{\mathbb{1}(\gamma < \theta)|\gamma \in \Gamma\}, \tag{4}$$

where $\mathbb{1}$ is the indicator function, $\Gamma$ is all the scaling factors of the network, and $\theta$ is the predefined pruning threshold of the pruning method. Actually, the pruning mask $\mathcal{M}$ consists of the unimportant-channel mask of each layers, *i.e.* $\mathcal{M} = \{\mathcal{M}^{(1)}, \mathcal{M}^{(2)}, ..., \mathcal{M}^{(L)}\}$, where $L$ is the layer count of the network. According to Equation 4, the pruning masks $\mathcal{M}_i$ is a binary vector consisting of 0 and 1.

**Algorithm 1** Algorithm Description of MaskSparsity

---

**Input:** training data: $\{\mathcal{X}, y\}$, pruning threshold $\theta$.

1: **Initialize**: pretrained model parameter $\mathcal{W} = \{\mathcal{W}_i, 0 \leq i \leq L\}$;

      **For:**$epoch = 1; epoch \leq epoch_{max}; epoch + +$

2: Update the model parameter $\mathcal{W}$ based on $\{\mathcal{X}, y\}$ and, using the global sparse regularization as in Equation 3; **Endfor:**

3: Obtain the pruning mask $\mathcal{M}$ by thresholding $\gamma$ with $\theta$;

4: **Reinitialize**: pretrained model parameter $\mathcal{W} = \{\mathcal{W}_i, 0 \leq i \leq L\}$;

      **For:**$epoch = 1; epoch \leq epoch_{max}; epoch + +$

5: Update the model parameter $\mathcal{W}$ based on $\{\mathcal{X}, y\}$ , the mask-guided sparse regularization as shown in Equation 5 with the mask $\mathcal{M}$; **Endfor:**

6: Obtain the compact model $\mathcal{W}^*$ from $\mathcal{W}$;

7: Finetune the compact model $\mathcal{W}^*$;

**Output:** The compact model and its parameters $\mathcal{W}^*$.

---

As discussed above, the over-regularization of the important channels limits the network capacity. Therefore, in this paper, we design a fine-grained sparse training strategy to alleviate the damage of the sparse regularization loss on important channels. Specifically, we propose to apply the sparse regularization only on the unimportant channels. Based on Equation 3, we can describe our MaskSparsity method as the Equation 5:

$$L_{\text{Sparsity}}(\mathcal{I}^0, y, \mathcal{W}) = L(f(\mathcal{I}^0, \mathcal{W}), y) + \lambda \cdot \sum_{l=1}^{L} \sum_{i=1}^{N_l} \mathcal{M}_i^{(l)} ||\gamma_i^{(l)}||_g ,$$
(5)

where $\mathcal{M}$ denotes the binary mask, indicating the unimportant channels of the whole network. For the important channels, the values in the channel mask are 0. Therefore, these channels are not affected by the sparse regularization and are trained as normal.

## 4. Experiments

### 4.1. Experimental Settings

**Datasets**. To demonstrate the effectiveness of Mask Sparsity in reducing model complexity, we evaluate MaskSparsity on both small and large datasets, *i.e.*, CIFAR-10 [16] ,and ILSVRC-2012 [30]. The CIFAR-10 dataset consists of natural images of 10 classes with resolution $32 \times 32$ and the train and test sets contain 50,000 and 10,000 images respectively. The ImageNet dataset consists natural images of 1000 classes with resolution $224 \times 224$ and the train and test sets contain 1.2 million and 50,000 images respectively. We experiment with ResNet-50 [9] on ILSVRC-2012, and experiment with ResNet-56 [8] and ResNet-110

[8] on CIFAR-10.

**Codebase and Baseline**. We directly deploy our algorithm on two popular codebases in github[1,2] for the experiments on cifar-10, cifar-100, and ILSVRC-2012. We hardly ever change any of the origin codes, except for adding the codes of our algorithm. Due to the difference in the training strategy of the baseline with other methods, we have the higher baseline accuracy on the evaluation datasets. It should be pointed that a higher baseline makes it difficult for the pruning algorithms to keep the accuracy after pruning.

**Evaluation Protocols**. We use the number of parameters and the FLOPS [8] to evaluate the complexity of the networks. To evaluate the accuracy, we use top-1 and top-5 score of full-size models and pruned models on ILSVRC-2012 and top-1 score only on CIFAR-10.

**Training and Pruning setting**. All the training-related hyper-parameters follow the two above-mentioned GitHub repositories. We use the same hyper-parameters during the normal training stage, the two sparsity training stages, and the finetuning stage in Figure 2. Specifically, on CIFAR-10, we train models for 200 epochs with a batch size of 128, a weight decay of 0.0005, a Nesterov momentum of 0.9 without dampening in every stage, and an initial learning rate of 0.1 which is divided by 5 at epochs 60, 120 and 160 on four NVIDIA GTX 1080Ti GPUs; On ILSVRC-2012, we train models for 100 epochs with a batch size of 256, a weight decay of 0.0001, a Nesterov momentum of 0.9 with dampening in every stage, and an initial learning rate of 0.1 which is divided by 10 at epochs 30, 60 and 90 on eight GPUs.

We set $\lambda$ as $2e^{-4}$ and $5e^{-4}$ separately for the global sparsity training stage and the mask sparsity training stage. We only manually set them without too much tuning. We think the former should be set lower since it affects all channels of each layer. Moreover, in the two sparsity training stages, we reinitialize the network with a normally-trained model.

For the pruning mask generation step after the global sparse regularization stage, we use a threshold of $1e^{-2}$. This thresholding step is to generate a sparse mask for the mask sparse regularization stage. Also, we use this sparse mask as our final pruning mask after the mask sparse regularization. In this way, we perform the *pruning-aware* sparse regularization on the network.

After pruning, we fine-tune the pruned models with an initial learning rate of 0.001, and keep other parameter settings the same as the previous step, on both the datasets.

---

[1]https://github.com/weiaicunzai/pytorch-cifar100
[2]https://github.com/facebookresearch/pycls

Table 1. Evaluation results using ResNet-50 on ILSVRC-2012.

| Method | Base Top-1(%) | Base Top-5(%) | Pruned Top-1(%) | Pruned Top-5(%) | Top-1 ↓ (%) | Top-5 ↓ (%) | FLOPs ↓ (%) |
|--------|-----|-----|-----|-----|-----|-----|-----|
| NS [26] | 75.04 | - | 69.60 | - | 5.44 | - | 50.51 |
| OT [37] | 75.04 | - | 70.40 | - | 4.64 | - | 52.88 |
| SFP [11] | 76.15 | 92.87 | 74.61 | 92.06 | 1.54 | 0.81 | 41.8 |
| GAL-0.5 [23] | 76.15 | 92.87 | 71.95 | 90.94 | 4.20 | 1.93 | 43.03 |
| HRank [22] | 76.15 | 92.87 | 74.98 | 92.33 | 1.17 | 0.54 | 43.76 |
| Hinge [20] | - | - | 74.7 | - | - | | 46.55 |
| HP [35] | 76.01 | 92.93 | 74.87 | 92.43 | 1.14 | 0.50 | 50 |
| MetaPruning [27] | 76.6 | - | 75.4 | - | 1.2 | - | 51.10 |
| Autopr [29] | 76.15 | 92.87 | 74.76 | 92.15 | 1.39 | 0.72 | 51.21 |
| FPGM [13] | 76.15 | 92.87 | 74.83 | 92.32 | 1.32 | 0.55 | 53.5 |
| DCP [41] | 76.01 | 92.93 | 74.95 | 92.32 | 1.06 | 0.61 | 55.76 |
| ThiNet [28] | 75.30 | 92.20 | 72.03 | 90.99 | 3.27 | 1.21 | 55.83 |
| EagleEye*[1] [17] | 77.21 | 93.68 | 76.37 | 92.89 | 0.84 | 0.79 | 50 |
| **MaskSparsity (ours)** | **76.44** | **93.22** | **75.68** | **92.78** | **0.76** | **0.44** | **51.07** |
| HRank [22] | 76.15 | 92.87 | 71.98 | 91.01 | 4.17 | 1.86 | 62.10 |

[1] The baseline of EagleEye* is obtained by evaluating the weight provided by the authors.

## 4.2. Results and Analysis

### 4.2.1 Results on ILSVRC-2012

As shown in Table 1, our proposed MaskSparsity achieves the new state-of-the-art. NS [26] is the baseline method of MaskSparsity, which adopts the global sparse regularization on the scaling factors of BN layers. OT [37] is the improved NS method, which sets an optimal threshold for each layer. It can be seen that the MaskSparsity outperforms them significantly in the respect of accuracy drop (Top-1 ↓ and Top-5 ↓ in Table 1) under roughly the same level of FLOPS drop. This shows that with the fine-grained sparse regularization, we avoid the bad effect of the sparse regularization on the unpruned channels. Moreover, as shown in the ablation study 4.3, with MaskSparsity, the pruning threshold on the scaling factors is easier to set.

In Table 1, it can be seen that we also outperform the non-sparsity-training-based methods under the same FLOPS decrease rate, e.g., FPGM [13] (53.5% FLOPS reduced), DCP [41] (55.76% FLOPS reduced), MetaPruning [27] (51.10% FLOPS reduced), and EagleEye [17] (50% FLOPS reduced). While the FLOPS reduction of MaskSparsity is less than HRank (53.76% vs 62.1%), the accuracy of the pruned model is much higher than that of HRank (0.93 vs 4.17 in Top-1 accuracy drop). This shows MaskSparsity's superiority over the previous state-of-the-art methods.

### 4.2.2 Results on CIFAR-10

Table 2 shows the experimental results of ResNet-56 on CIFAR-10. On this small dataset, MaskSparsity also achieves the state-of-the-art performance. Under similar FLOPs reduction with FPGM [13] and Hinge [20],

Table 2. Evaluation results using ResNet-56 on CIFAR-10.

| Method | Base Top-1(%) | Pruned Top-1 (%) | Top-1 ↓ (%) | FLOPs ↓ (%) |
|--------|-----|-----|-----|-----|
| NISP [39] | - | - | 0.03 | 42.6 |
| Hinge [20] | 93.69 | 92.95 | 0.74 | 50 |
| AMC [12] | 92.8 | 91.9 | 0.9 | 50 |
| LeGR [3] | 93.9 | 93.7 | 0.2 | 52 |
| FPGM [13] | 93.59 | 93.26 | 0.33 | 52.6 |
| LFPC [10] | 93.59 | 93.24 | 0.35 | 52.9 |
| **MaskSparsity (ours)** | **94.50** | **94.19** | **0.31** | **54.88** |
| GAL-0.8 [23] | 93.26 | 90.36 | 2.9 | 60.2 |
| HRank [22] | 93.26 | 90.72 | 2.54 | 74.1 |

MaskSparsity achieves 0.31% Top-1 accuracy drop with ResNet-56, which is slightly better than FPGM [13] (0.31% vs 0.33%) and Hinge [20] (0.31% vs 0.74%).

Table 3 shows the experimental results of another network, ResNet-110 on CIFAR-10. With this deeper network, MaskSparity achieves a better performance. As shown in Table 3, our MaskSparsity outperforms the other state-of-the-art methods, like HRank [22] (at 58.2% FLOPS reduction), under roughly the same ratio of FLOPS reduction. MaskSparsity has roughly the same accuracy increase with FPGM [13] (0.02 vs 0.06 ), but MaskSparsity has a larger FLOPS reduction than these two methods(63.03% vs 52.3% and 60.89%).

Table 4 shows the experimental results of VGG-16, which is a straight network structure that is different from ResNet. We compare MaskSparsity with NetSlim (NS) [26], FPGM [13], and PFEC [19]. It shows that we outperforms NS [26] and PFEC [19] on both accuracy and FLOPS reduction. Compared with FPGM [13], we are 0.04 % less

Table 3. Evaluation results using ResNet-110 on CIFAR-10.

| Method | Base Top-1(%) | Pruned Top-1(%) | Top-1 ↓ (%) | FLOPs↓ (%) |
|---|---|---|---|---|
| Li et al. [19] | 93.53 | 93.30 | 0.23 | 38.60 |
| SFP [11] | 93.68 | 93.86 | -0.18 | 40.8 |
| NISP-110 [39] | - | - | 0.18 | 43.78 |
| GAL-0.5 [23] | 93.50 | 92.74 | 0.76 | 48.5 |
| FPGM [13] | 93.68 | 93.74 | -0.06 | 52.3 |
| HRank [22] | 93.50 | 93.36 | 0.14 | 58.2 |
| LFPC [10] | 93.68 | 93.07 | 0.61 | 60.3 |
| **MaskSparsity (ours)** | **94.70** | **94.72** | **-0.02** | **63.03** |
| HRank [22] | 93.50 | 92.65 | 0.85 | 68.6 |
| SASL [31] | 93.83 | 93.80 | 0.03 | 70.2 |

Table 4. Evaluation results using VGG-16 on CIFAR-10. FT: fine-tuning.

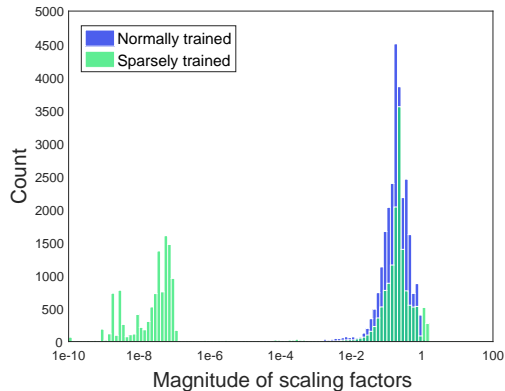| Method | Base Top-1(%) | before FT(%) | FT (%) | Top-1 ↓ (%) | FLOPs↓ (%) |
|---|---|---|---|---|---|
| PFEC [19] | 93.58 | 77.45 | 93.28 | 0.3 | 34.2 |
| FPGM [13] | 93.58 | 80.38 | 94.00 | -0.42 | 34.2 |
| NS [26] | 93.66 | - | 93.80 | -0.14 | 51 |
| **MaskSparity(ours)** | **93.86** | **94.16** | **94.24** | **-0.38** | **52.21** |



Figure 4. Distribution of scaling factors of ResNet50 on ILSVRC-2012 before and after the MaskSparsity's sparsity training.

than FPGM on the increase of the accuracy, which is very minor. However, we decrease 18.01% more FLOPS than FPGM. Therefore, we also outperform FPGM in general pruning performance. Moreover, we also compare the accuracy drop of the pruned model without the fine-tuning stage. It can be seen in the third column of Table 4 that the MaskSparsity suffers a weaker accuracy drop than the other methods.

### 4.3. Ablation Study

**Visualization of the distribution of the scaling factors after using MaskSparsity.** In Section 3.3 and Figure 3, we show that the distribution of scaling factors meets the over-regularization problem that might damage the pruning performance. In Figure 4 we draw the distribution of the same set of scaling factors after the mask-guided sparse regular-
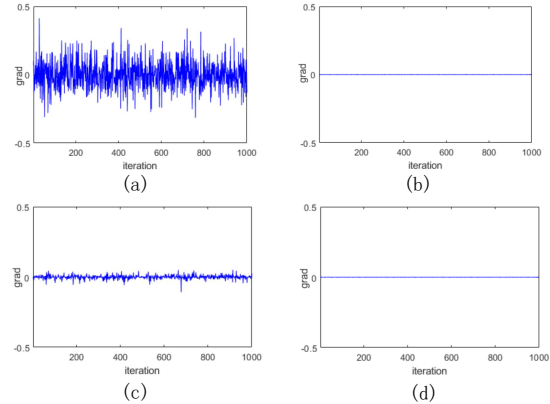


Figure 5. The gradients of the scaling factors of some certain channel at the end of the sparsity training stage. Both important and unimportant channels are visualized. (a) The important channel, global sparsity. (b) The unimportant channel, global sparsity. (c) The important channel, MaskSparsity. (d) The unimportant channel, MaskSparsity.

ization and compare it with that of the pre-trained networks. It can be seen that the two problems are alleviated significantly. The right peak of the sparsity trained network does not move towards 0 like in Figure 3. Moreover, the two peaks are well distinguishable, with almost no in-between bars in the middle area. This demonstrates the effectiveness of the fine-grained sparse regularization of MaskSparsity, which would benefit the pruning performance.

**The convergence analysis by visualizing the gradients.** To validate the statement that the sparse regularization on unpruned channels restricts the network's capacity, we visualize the gradients of important and unimportant channels after the sparsity training using global and fine-grained sparse regularization, respectively. The result is shown in Figure 5. The gradients are collected by continuing training for 1,000 iterations from the end of the sparsity training stage of ResNet-50 on ILSVRC-2012. In Figure 5 (a), it can be seen that the gradient norm of the important channels is still large for the important channels with the global sparse regularization, while Figure 5 (b) shows the important channel has a small gradient norm with our fine-grained MaskSparsity sparse regularization. Figure 5 (b) and (d) shows that both the gradients of unimportant channels with the two kinds of sparse regularization methods are almost the same. According to these figures, although the network's accuracy and sparsity have converged, the global sparse regularization leads to a larger gradient on the important channels. We infer that the large gradient prevents the network converges to a better local minimum.

**Comparing fine-tuning and train-from-scratch.** In Table 5, we list the network model's accuracy and computational complexity at different pruning stages. Moreover, we also list the result that trains the pruned model from scratch

Table 5. The stage-wise performance in the case of pruning ResNet-56 on CIFAR-10.

| Model state | Top-1(%) | FLOPs | Parameters |
|---|---|---|---|
| Normally trained | 94.50 | 126M | 853K |
| MaskSparsity trained | 94.02 | 126M | 853K |
| Pruned | 92.67 | 57M | 419K |
| Finetune | 94.19 | 57M | 419K |
| Train from Scratch | 93.60 | 57M | 419K |

Table 6. Evaluation results of MaskSparsity using the pruning mask of EagleEye.

| Model state | Top-1(%) | Top-5(%) | FLOPs↓(%) |
|---|---|---|---|
| Unpruned Resnet-50 | 77.21 | 93.68 | - |
| EagleEye [17] | 76.37 | 92.89 | 50.0 |
| MaskSparsity + EagleEye | 76.69 | 93.22 | 50.0 |

Table 7. Evaluation results of MaskSparsity based on uniform pruning mask of ResNet-50 on ImageNet.

| Model state | Top-1 (%) | Top-5 (%) | FLOPs↓ (%) |
|---|---|---|---|
| Unpruned Resnet-50 | 76.44 | 93.22 | - |
| Direct pruning with Uniform Mask | 74.23 | 92.28 | 53.46 |
| MaskSparsity with Uniform Mask | 75.62 | 92.68 | 53.46 |

Table 8. Evaluation results of MaskSparsity with different regularizations of ResNet-56 on CIFAR-10.

| Model state | Top-1(%) | Top-5(%) | FLOPs↓(%) |
|---|---|---|---|
| Unpruned Resnet-56 | 94.50 | 99.79 | - |
| MaskSparsity with L1 | 94.19 | 99.81 | 54.88 |
| MaskSparsity with L2 | 93.99 | 99.8 | 54.88 |

Table 9. Evaluation results of MaskSparsity on an YOLOv5s-based face detector on WiderFace.

| Model state | mAP[Easy] | FLOPS | Param |
|---|---|---|---|
| YOLOv5s | 92.38% | 4.1G | 3.5M |
| YOLOv5s+MaskSparsity | 91.86% | 2.0G | 1.6M |

without exploiting its weights. It can be seen that the train-from-scratch result is lower than the fine-tuning result. We think this demonstrates the effectiveness of the fine-grained sparsely-trained pre-trained weights.

**The performance on the pruning mask generators.** As discussed above, the above pruning process consists of two key elements, i.e. identifying the unimportant channels and pushing them to 0 by sparse regularization. This paper uses global sparsity training to generate the pruning mask. To show the MaskSparsity's generalization on other pruning mask generators, we apply it to two other pruning methods and show the superiority on the performance. Except for the pruning mask generating method, the other stage is the same as the pipeline in Figure 2.

Firstly, we directly use the codebase of EagleEye [17] and use the pruning mask after its searching process to conduct our mask sparse regularization stage. We reuse the hyper-parameters of EagleEye's finetuning stage in our sparse regularization. Table 6 shows the experimental results. Under the same pruning mask after pruning, the top-1 accuracy increased by 0.32% over the original EagleEye. This shows the scalability of MaskSparsity on state-of-the-art methods.

Secondly, we try the naive pruning method that directly pruning the same portion of channels of each layer. We call this naive pruning method as *Uniform Pruning*. Experimental results on ResNet-50 are shown in Table 7. It can be seen that MaskSparsity improves this naive pruning method to SOTA-level performance. This demonstrates the effectiveness of MaskSparsity on pushing the unimportant channels to near 0 without too much damage on the important channels.

**Masksparsity with different Regularization.** In this paper, we mainly use the L1 regularization to generate the network sparsity. To show the compatibility of MaskSparsity with other specific forms of regularizations, we replace the L1 regularization with L2 regularization in the mask sparse training stage. We conduct this ablation study using ResNet-56 on CIFAR-10. The experimental results are shown in Table 8. It can be seen that there is little difference in accuracy between the result of L1 regularization in MaskSparsity stage and L2 regularization. The accuracy of L1 regularization in MaskSparsity stage is 0.2 points higher than L2 regularization.

### 4.4. Application on Other Tasks

To validate the generalization ability, We apply the method to two different object detection tasks. The first is the face detection based on YOLOv5[3] evaluated on Wider-Face [36]. The second is the car detection based on Faster-RCNN-FPN [24] evaluated on PASCAL VOC. The results are listed in Table 9 and Table 10. For Faster-RCNN-FPN, we only prune the backbone part and report the backbone's FLOPS and parameters. From these experimental results, the pruned models of both tasks maintain roughly the same level of accuracy. It can be concluded that MaskSparsity applies to other tasks.

### 5. Conclusion

In this paper, to solve the problem that existing sparsity-training-based methods over-regularize the important chan-

---

[3]https://github.com/ultralytics/yolov5

Table 10. Evaluation results of MaskSparsity on an FPN-based car detector on PASCAL VOC. The input size is 1000×600 and the backbone is ResNet-50.

| Model state | mAP | FLOPS | Param |
|---|---|---|---|
| Faster-RCNN-FPN | 89.7% | 49.95G | 23.5M |
| Faster-RCNN-FPN+MaskSparsity | 89.3% | 23.73G | 11.4M |

nels, we design a pruning-aware sparse training method, named as MaskSparsity. MaskSparsity only applies the sparse regularization on the unimportant channels which are to be pruned. Therefore, MaskSparsity can minimize the negative impact of the sparse regularization on the important channels. The method is effective and efficient. The experimental results show that it outperforms the other sparsity-training-based pruning methods and achieves the state-of-the-art on the benchmarks. In the future, we plan to work on how to obtain better pruning masks.

# References

[1] Jose M Alvarez and Mathieu Salzmann. Learning the number of neurons in deep networks. In *Advances in Neural Information Processing Systems*, pages 2270–2278, 2016. 1, 2

[2] Jose M Alvarez and Mathieu Salzmann. Compression-aware training of deep networks. In *Advances in Neural Information Processing Systems*, pages 856–867, 2017. 3

[3] Ting-Wu Chin, Ruizhou Ding, C. Zhang, and Diana Marculescu. Towards efficient model compression via learned global ranking. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1515–1525, 2020. 6

[4] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1

[5] P. T. Fletcher, Suresh Venkatasubramanian, and S. Joshi. Robust statistics on riemannian manifolds via the geometric median. *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, 2008. 3

[6] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017. 4

[7] Song Han, Jeff Pool, John Tran, and W. Dally. Learning both weights and connections for efficient neural network. *ArXiv*, abs/1506.02626, 2015. 4

[8] 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. 1, 5

[9] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. *ArXiv*, abs/1603.05027, 2016. 5

[10] Yang He, Yuhang Ding, Ping Liu, Linchao Zhu, Hanwang Zhang, and Yi Yang. Learning filter pruning criteria for deep convolutional neural networks acceleration. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 6, 7

[11] Yang He, Guoliang Kang, Xuanyi Dong, Yanwei Fu, and Yi Yang. Soft filter pruning for accelerating deep convolutional neural networks. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence*, pages 2234–2240, 2018. 6, 7

[12] Yihui He, Ji Lin, Zhijian Liu, Hanrui Wang, Li-Jia Li, and Song Han. Amc: Automl for model compression and acceleration on mobile devices. In *European Conference on Computer Vision*, pages 815–832. Springer, 2018. 6

[13] Yang He, Ping Liu, Ziwei Wang, Zhilan Hu, and Yi Yang. Filter pruning via geometric median for deep convolutional neural networks acceleration. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 4340–4349. Computer Vision Foundation / IEEE, 2019. 3, 6, 7

[14] Zehao Huang and Naiyan Wang. Data-driven sparse structure selection for deep neural networks. In *Proceedings of the European conference on computer vision (ECCV)*, pages 304–320, 2018. 1, 3

[15] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. 3

[16] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, 2009. 5

[17] Bailin Li, Bowen Wu, Jiang Su, Guangrun Wang, and Liang Lin. Eagleeye: Fast sub-net evaluation for efficient neural network pruning. In *ECCV*, 2020. 6, 8

[18] Hao Li, Asim Kadav, Igor Durdanovic, H. Samet, and H. Graf. Pruning filters for efficient convnets. *ArXiv*, abs/1608.08710, 2017. 3

[19] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. In *5th International Conference on Learning Representations*, 2017. 6, 7

[20] Yawei Li, Shuhang Gu, C. Mayer, L. Gool, and R. Timofte. Group sparsity: The hinge between filter pruning and decomposition for network compression. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8015–8024, 2020. 3, 6

[21] Mingbao Lin, Rongrong Ji, Yan Wang, Yichen Zhang, Baochang Zhang, Yonghong Tian, and Ling Shao. Hrank: Filter pruning using high-rank feature map. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1529–1538, 2020. 3

[22] Mingbao Lin, Rongrong Ji, Yan Wang, Yichen Zhang, Baochang Zhang, Yonghong Tian, and Ling Shao. Hrank: Filter pruning using high-rank feature map. *CoRR*, abs/2002.10179, 2020. 6, 7

[23] Shaohui Lin, Rongrong Ji, Chenqian Yan, Baochang Zhang, Liujuan Cao, Qixiang Ye, Feiyue Huang, and David S. Doermann. Towards optimal structured CNN pruning via generative adversarial learning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 2790–2799. Computer Vision Foundation / IEEE, 2019. 6, 7

[24] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. 8

[25] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755, 2014. 1

[26] Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang. Learning efficient convolutional networks through network slimming. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2736–2744, 2017. 1, 3, 6, 7

[27] Zechun Liu, Haoyuan Mu, Xiangyu Zhang, Zichao Guo, Xin Yang, Kwang-Ting Cheng, and Jian Sun. Metapruning: Meta learning for automatic neural network channel pruning. In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*, pages 3295–3304. IEEE, 2019. 6

[28] Jian-Hao Luo, Hao Zhang, Hong-Yu Zhou, Chen-Wei Xie, Jianxin Wu, and Weiyao Lin. Thinet: Pruning CNN filters for a thinner net. *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(10):2525–2538, 2019. 6

[29] Jian-Hao Luo and Jianxin Wu. Autopruner: An end-to-end trainable filter pruning method for efficient deep model inference. *arXiv preprint arXiv:1805.08941*, 2018. 6

[30] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. 1, 5

[31] Jun Shi, Jianfeng Xu, Kazuyuki Tasaka, and Zhibo Chen. SASL: saliency-adaptive sparsity learning for neural network acceleration. *CoRR*, abs/2003.05891, 2020. 7

[32] Suraj Srinivas, Akshayvarun Subramanya, and R. Venkatesh Babu. Training sparse neural networks. *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 455–462, 2017. 3

[33] Enzo Tartaglione, S. Lepsøy, Attilio Fiandrotti, and G. Francini. Learning sparse neural networks via sensitivity-driven regularization. In *NeurIPS*, 2018. 4

[34] Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity in deep neural networks. In *Advances in neural information processing systems*, pages 2074–2082, 2016. 1, 2

[35] Xiaofan Xu, Mi Sun Park, and Cormac Brick. Hybrid pruning: Thinner sparse networks for fast inference on edge devices. *arXiv preprint arXiv:1811.00482*, 2018. 6

[36] Shuo Yang, Ping Luo, Chen-Change Loy, and Xiaoou Tang. Wider face: A face detection benchmark. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5525–5533, 2016. 8

[37] Yun Ye, Ganmei You, J. Fwu, Xia Zhu, Q. Yang, and Y. Zhu. Channel pruning via optimal thresholding. *ArXiv*, abs/2003.04566, 2020. 3, 6

[38] Yun Ye, Ganmei You, Jong-Kae Fwu, Xia Zhu, Qing Yang, and Yuan Zhu. Channel pruning via optimal thresholding. *arXiv preprint arXiv:2003.04566*, 2020. 4

[39] Ruichi Yu, Ang Li, Chun-Fu Chen, Jui-Hsin Lai, Vlad I Morariu, Xintong Han, Mingfei Gao, Ching-Yung Lin, and Larry S Davis. Nisp: Pruning networks using neuron importance score propagation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9194–9203, 2018. 6, 7

[40] Kai Zhao, Xinyu Zhang, Q. Han, and Ming-Ming Cheng. Dependency aware filter pruning. *ArXiv*, abs/2005.02634, 2020. 3

[41] Zhuangwei Zhuang, Mingkui Tan, Bohan Zhuang, Jing Liu, Yong Guo, Qingyao Wu, Junzhou Huang, and Jinhui Zhu. Discrimination-aware channel pruning for deep neural networks. In *Advances in Neural Information Processing Systems*, pages 883–894, 2018. 6