

# Softer Pruning, Incremental Regularization

Linhang Cai<sup>\*†</sup>, Zhulin An<sup>\*‡</sup>, Chuanguang Yang<sup>\*†</sup> and Yongjun Xu<sup>\*</sup>

<sup>\*</sup>Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China

<sup>†</sup>University of Chinese Academy of Sciences, Beijing, China

Email: {cailinhang19g, anzhulin, yangchuanguang, xyj}@ict.ac.cn

**Abstract**—Network pruning is widely used to compress Deep Neural Networks (DNNs). The Soft Filter Pruning (SFP) method zeroizes the pruned filters during training while updating them in the next training epoch. Thus the trained information of the pruned filters is completely dropped. To utilize the trained pruned filters, we proposed a Softer Filter Pruning (SRFP) method and its variant, Asymptotic Softer Filter Pruning (ASRFP), simply decaying the pruned weights with a monotonic decreasing parameter. Our methods perform well across various networks, datasets and pruning rates, also transferable to weight pruning. On ILSVRC-2012, ASRFP prunes 40% of the parameters on ResNet-34 with 1.63% top-1 and 0.68% top-5 accuracy improvement. In theory, SRFP and ASRFP are an incremental regularization of the pruned filters. Besides, We note that SRFP and ASRFP pursue better results while slowing down the speed of convergence.

## I. INTRODUCTION

Currently, deep Convolutional Neural Networks (CNNs) have shown extraordinary performance in various tasks, e.g., image classification [1], [2], target detection [3], semantic segmentation [4]. However, large burdens of DNN model size, limited run-time memory and huge numbers of Floating Point Operations (FLOPs) [5] hinder the deployment of DNN models in mobile devices. Thus, it matters to compress the DNN models to improve the computational efficiency. Prevalent methods of model compression and acceleration include low rank approximation [6] and network pruning [7]. Among them, filter pruning can reduce both the model size and the computational cost, thus, having been a hot research topic.

In terms of filter pruning, Soft Filter Pruning (SFP) maintains the capacity of the DNNs while pruning [8]. SFP zeroizes the filters chosen to be pruned and updates them in the next training epoch, as shown in Figure 1. However, traditional hard filter pruning (HFP) method would not use those pruned filters any more.

A drawback of SFP is that there is a severe accuracy drop after pruning in case of large pruning rates. Asymptotic Soft Filter Pruning (ASFP) is a variant of SFP to stabilize the training and pruning process [9]. ASFP gradually increases the pruning rate towards the objective pruning rate to reduce the information loss caused by setting pruned filters to zeros while pruning.

To utilize the trained pruned filters, we proposed a Softer Filter Pruning (SRFP) method and its asymptotic version, Asymptotic Softer Filter Pruning (ASRFP), simply decaying

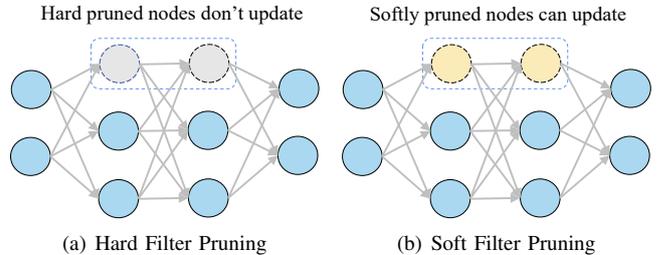


Fig. 1. Comparison of hard filter pruning and soft filter pruning. The gray nodes removed by hard filter pruning could not update in the next training epoch, while the yellow nodes pruned by soft filter pruning method can still update.

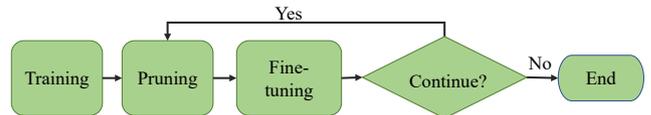


Fig. 2. A typical three-step pruning pipeline composed of three phases: training, pruning and fine-tuning. Our SRFP or ASRFP is used in the pruning phase to remove those filters chosen to be pruned smoothly using weights that gradually decay to zero, while the conventional pruning operation simply sets pruned filters to zeros.

the pruned weights with a monotonic decreasing parameter. SRFP is indeed an incremental  $\ell_2$ -norm regularization of the pruned filters.

A typical three-step pruning pipeline consists of three phases: training, pruning and fine-tuning, as shown in Figure 2. Our SRFP or ASRFP is used in the pruning phase to remove those filters chosen to be pruned smoothly using weights that gradually decay to zero.

While the conventional pruning operation is to simply set parameters chosen to be pruned to zeros, our method ensures that the pruned filters are removed smoothly using weights that gradually decay to zero, so that we can better preserve the trained information in those filters. Our method, based on SFP, only needs tiny extra computations and hyperparameters, so it is easy and effective to use our method to prune a model.

With the same goal of finding minimal nets, our method is very different from the weight decay regularization [10]. Actually, they are used in different phases. While weight decay is used in the training or fine-tuning phase to increase the sparsity of the networks or to avoid over-fitting, our method is used in the pruning phase to punish those pruned parameters in a soft manner. In short, while weight decay is used to constrain all weights in a network in the training or fine-tuning

<sup>‡</sup>Zhulin An is the corresponding author.

phase, our method only constrains those weights chosen to be pruned in the pruning phase. Our method will not affect those parameters that are not chosen to be pruned. In fact, in all of our experiments, weight decay is widely used in training and fine-tuning to avoid over-fitting.

Our contributions are as follows: (1) Gradually decaying the importance of the pruned filters, our methods perform well across various networks, datasets and pruning rates, also transferable to weight pruning. (2) The gradually decaying manner of SRFP and ASRFP is actually doing incremental regularization on those pruned weights, reserving the pruned weights at start, gradually forcing them towards zeros. (3) We note that SRFP, ASRFP and ASFP pursue better results while slowing down the speed of convergence.

## II. RELATED WORKS

Previous attempts on deep neural network compression mainly include *low rank matrix factorization*, *fast convolution* and *pruning*.

Among them, *low rank matrix factorization* is a tensor low rank expansion technique to reduce the number of parameters or speedup deep neural networks [6]. However, these methods reveal relatively small speedups on small convolutional kernels such as  $3 \times 3$  and  $1 \times 1$ , which are widely used in prevalent CNNs like ResNets with bottleneck structures [1]. Compression-aware training [11] uses a regularizer to push the parameter matrix of each layer to be a low rank matrix as close as possible.

*Fast Convolution* finds various efficient convolutional filters to design the efficient architecture, including Groupwise Convolution [12], Depthwise Separable Convolution [13], and Heterogeneous Kernel-Based Convolution (HetConv) [14].

*Pruning* techniques focus on reducing the network complexity by removing unimportant neural nodes or connections [7], [15]. Many studies on network pruning calculate the importance of the filters or connections and prunes them based on some criteria, and then finetune the pruned network to avoid severe accuracy drop. Some studies explore the automated determination of the threshold values for pruning [16].

**Weight Pruning.** Magnitude-based weight pruning methods are computationally efficient, compressing networks by deleting unimportant weights in an unstructured manner. Iterative weight pruning is a three-step method to discard small weights whose magnitude below the threshold [17]. Dynamic network surgery properly incorporates both the pruning and splicing operations for model compression to avoid incorrect pruning, instead of alternately pruning and retraining [18].

**Filter Pruning.** Filter pruning removes redundant filters (channels) or feature maps. Prevalent metrics to evaluate the filter importance include  $\ell_1$ -norm,  $\ell_2$ -norm, scaling factors and feature redundancy [19], [5], [20]. Gate Decorator [21] multiplies the output of a CNN module by channel-wise scaling factors, using Taylor expansion to evaluate the impact on the loss function owing to setting the scaling factor of each filter to zero. A common drawback of these pruning techniques is that the network capacity is decreased after pruning.

AutoPruner [22] combines the pruning phase and the fine-tuning phase of the typical three-step pruning pipeline into one single end-to-end trainable system to find unimportant filters automatically during training, thus increasing the computational and memory consumptions.

SFP and ASFP method set the pruned filters to zeros during training while updating them in the next training epoch to maintain the representative capacity.

During the first several pruning epochs of SFP, there is an obvious accuracy loss in the test set. ASFP gradually increases the pruning rate towards the aimed pruning rate to remedy this issue. To better utilize the trained pruned filters, we proposed a Softer Filter Pruning (SRFP) method and its asymptotic version, Asymptotic Softer Filter Pruning (ASRFP), simply decaying the pruned weights with a monotonic decreasing parameter  $\alpha$ . SRFP is an incremental regularization of the pruned filters. Our approach works well without much computational and memory consumptions.

## III. OUR METHOD

### A. Formulation

Consider the convolutional kernel  $W_i \in \mathbb{R}^{n \times m \times s \times s}$  in the  $i$ -th layer, where  $1 \leq i \leq L$  and  $L$  is the number of convolutional layers. Specifically,  $s$ ,  $m$  and  $n$  are the convolutional kernel size, the number of input channels and output channels respectively.

The input feature map  $I_i \in \mathbb{R}^{m \times h_i \times w_i}$  and the output feature map  $O_i \in \mathbb{R}^{n \times h_{i+1} \times w_{i+1}} = W_i * I_i$  is calculated by

$$O_{i,j} = W_{i,j} * I_i \quad \text{for } 1 \leq j \leq n, \quad (1)$$

where  $O_{i,j} \in \mathbb{R}^{h_{i+1} \times w_{i+1}}$  and  $W_{i,j} \in \mathbb{R}^{m \times s \times s}$  denote the  $j$ -th output channel and the  $j$ -th filter of the  $i$ -th layer respectively.

Suppose that the filter pruning rate for the  $i$ -th layer is  $P_i$ . Thus there are  $n \times P_i$  filters to be removed in the  $i$ -th layer, and the size of the pruned output feature map  $O_i$  would be  $n \times (1 - P_i) \times h_{i+1} \times w_{i+1}$ .

According to the SFP method and its variant ASFP, the pruned weights of the  $i$ -th layer are simply zeroized, which can be represented by

$$\hat{W}_{i,j} = W_{i,j} \odot M_{i,j} \quad \text{for } 1 \leq j \leq n, \quad (2)$$

where  $M_{i,j}$  is a Boolean matrix with the same shape as the filter  $W_{i,j}$  to denote whether the  $j$ -th filter is pruned or not in the  $i$ -th layer. Besides,  $\odot$  is a matrix pointwise multiplication operator. Exactly,  $M_{i,j} = 0$  if  $W_{i,j}$  is pruned. Otherwise, we let  $M_{i,j} = 1$  to denote that the filter  $W_{i,j}$  is not pruned.

### B. Motivation

Rewrite (2) equivalently as

$$\hat{W}_{i,j} = W_{i,j} \odot M_{i,j} + \alpha W_{i,j} \odot (1 - M_{i,j}) \quad \text{for } 1 \leq j \leq n, \quad (3)$$

where  $\alpha = 0$  in SFP as well as ASFP, and  $\alpha$  is the decaying rate for those pruned weights. We can replace  $\alpha$  with a decreasing nonzero number to better utilize the trained

information inside those pruned weights. In general, we set  $\alpha \in [0, 1]$ .

Notably, when  $\alpha \in (0, 1]$ , the trained knowledge of the pruned filters is not completely dropped which would be helpful for releasing the accuracy drop caused by the pruning phase and achieving better results in the next training epoch.

However, when  $\alpha > 0$ , the pruned filters are not zeroized so that the resulted pruned model is not compact. So we limit  $\alpha$  to gradually decay from the initial value  $\alpha_0$  where  $\alpha_0 \in [0, 1]$  towards zero as the training and pruning procedure goes on. We consider two kinds of decaying strategies, which are exponential decay and linear decay respectively.

The exponential decay strategy can be written as

$$\alpha_e(t) = \alpha_0 e^{\frac{-kt}{t_{max}-1}} \quad \text{for } 0 \leq t < t_{max}, \quad (4)$$

where  $t_{max}$  is the maximal number of training epochs and  $k$  is a coefficient to control the descent speed of  $\alpha$ . Then we introduce a constraint parameter  $\epsilon$  to obtain the value of  $k$ , claiming that

$$\alpha_e(t_{max} - 1) = \alpha_0 e^{-k} = \epsilon, \quad (5)$$

where  $\epsilon$  is infinitely close to zero. Thus  $k = \ln \frac{\alpha_0}{\epsilon}$  and  $\alpha_e(t)$  can be given by

$$\alpha_e(t) = \alpha_0 \left(\frac{\alpha_0}{\epsilon}\right)^{-\frac{t}{t_{max}-1}} \quad \text{for } 0 \leq t < t_{max}, \quad (6)$$

when  $\alpha_e(t)$  is approaching zero, we just set  $\alpha_e(t) = 0$  to obtain a really compact model. Similarly, the linear decay strategy can be given by

$$\alpha_l(t) = \alpha_0 \left(1 - \frac{t}{t_{max}-1}\right) \quad \text{for } 0 \leq t < t_{max}, \quad (7)$$

where  $\alpha_l(0) = \alpha_0$  and  $\alpha_l(t_{max} - 1) = 0$ .

### C. SofteR Filter Pruning (SRFP)

Based on the above motivation, we illustrate our SRFP method in Algorithm 1, where we prune  $P_i \times 100\%$  of the filters in the  $i$ -th convolutional layer according to the  $\ell_2$ -norm of all filters. For simplicity, we use the same pruning rate  $P_i$  for each convolutional layer to get rid of complicated hyper-parameter search.

At the beginning of the training and pruning phase, the pruned filters are decayed in a soft manner, especially when  $\alpha$  is close to 1, which means that we nearly maintain all the trained information inside the pruned filters. Thereby, we greatly avoid the sharp accuracy drop caused by pruning, achieving a better performance. As the phase goes on, we gradually push  $\alpha$  towards 0, making the training and pruning phase close to that of SFP method, in order to obtain an actually compact model.

Moreover, we can view our SRFP method from the angle of incremental regularization. Denote the pruned filters as

$$WP_{i,j} = W_{i,j} \odot (1 - M_{i,j}) \quad \text{for } 1 \leq j \leq n, \quad (8)$$

where  $1 - M_{i,j}$  is regarded as a mask to obtain those pruned filters of the  $i$ -th layer. Thus (3) equivalent to

$$\hat{W}_{i,j} = W_{i,j} \odot M_{i,j} + \alpha WP_{i,j} \quad \text{for } 1 \leq j \leq n, \quad (9)$$

where  $W_{i,j} \odot M_{i,j}$  is the completely maintained while the pruned portion  $WP_{i,j}$  is decayed by  $\alpha$ . We can split (9) into the following two steps:

$$\hat{W}P_{i,j} = \alpha WP_{i,j} \quad \text{for } 1 \leq j \leq n, \quad (10)$$

where  $\hat{W}P_{i,j}$  is  $WP_{i,j}$  decayed by  $\alpha$ , and then

$$\hat{W}_{i,j} = W_{i,j} \odot M_{i,j} + \hat{W}P_{i,j} \quad \text{for } 1 \leq j \leq n, \quad (11)$$

where  $\hat{W}_{i,j}$  is the initial value of the  $j$ -th filter in the  $i$ -th layer for the next training epoch.

Denote  $\lambda(t) = \alpha_0 - \alpha(t)$ , where  $\alpha(t)$  could be either exponential or linear decay. Thus we rewrite (10) equivalently as

$$\begin{aligned} \hat{W}P_{i,j} &= \alpha WP_{i,j} = (\alpha_0 - \lambda)WP_{i,j} = \alpha_0 WP_{i,j} - \lambda WP_{i,j} \\ &= \alpha_0 WP_{i,j} - \frac{\lambda}{2} \frac{\partial \|WP_{i,j}\|_2^2}{\partial WP_{i,j}} \quad \text{for } 1 \leq j \leq n. \end{aligned} \quad (12)$$

Consider the special case when  $\alpha_0 = 1$ . Our SRFP is indeed adding a  $\ell_2$ -norm regularization term  $\frac{\lambda}{2} \|WP_{i,j}\|_2^2$  to those pruned filters. Since  $\alpha(t)$  decreases from 1 to 0, then  $\lambda(t)$  increases from 0 to 1. So the regularization gradually strengthened.

Besides, as the asymptotic variant of SFP called ASFP gradually increases the pruning rate towards the final pruning rate, likewise, we also create an asymptotic version of SRFP named ASRFP that gradually increases the pruning rate.

---

#### Algorithm 1: SRFP Algorithm

---

**inputs :** training set:  $X$ , pruning rate:  $P_i$ , initial decay rate:  $\alpha$ , the model with parameters  $W = \{W_i, 0 \leq i \leq L\}$ .  
**output :** The pruned model with parameters  $W^* = W^{t_{max}}$   
Initialize the model parameter  $W^0$   
**for**  $t = 0, \dots, t_{max} - 1$  **do**  
    Decrease weight decay rate  $\alpha$   
    Train model parameters  $\hat{W}^{t+1}$  based on data set  $X$  and  $W^t$   
    **for**  $i = 1, \dots, L$  **do**  
        Compute the  $\ell_2$ -norm of each filter  $\|\hat{W}_{i,j}^{t+1}\|_2, 1 \leq j \leq n$   
        Select  $n \times P_i$  filters with minimal  $\ell_2$ -norm values  
        Decay the parameters of chosen filters with  $\alpha$   
    Get the softly pruned model parameters  $W^{t+1}$  based on  $\hat{W}^{t+1}$   
Get the pruned model with final parameters  $W^* = W^{t_{max}}$

---

## IV. EXPERIMENTAL RESULTS

### A. Setup

We evaluate our approaches on the CIFAR-10 [25], and ILSVRC-2012 [26]. CIFAR-10 includes 60,000 RGB images of size  $32 \times 32$  pixels, divided into 10 classes. ILSVRC-2012 consists of 1.28 million training images and 50k validation images drawn from 1,000 categories. We focus on pruning the prevalently utilized ResNet, following the common experimental setup in ThiNet [27], CP [23].

On CIFAR-10, we follow the parameter scheme and the training configuration in [28]. On ILSVRC-2012, we follow

TABLE I  
COMPARISON OF PRUNING RESNETS ON CIFAR-10

Depth	Method	Pre-trained?	Baseline Accu. (%)	Accelerated Accu. (%)	Accu. Drop (%)	FLOPs	Pruned FLOPs(%)
56	PFEC [19]	×	93.04	91.31	1.75	9.09E7	27.6
	CP [23]	×	92.80	90.90	1.90	-	50.0
	SFP(20%)	×	93.66 ± 0.28	93.26 ± 0.20	0.40	8.98E7	28.4
	ASFP(20%)	×	93.66 ± 0.28	93.26 ± 0.21	0.40	8.98E7	28.4
	Our1(20%)	×	93.66 ± 0.28	<b>93.33 ± 0.43</b>	<b>0.33</b>	8.98E7	28.4
	Our2(20%)	×	93.66 ± 0.28	<b>92.92 ± 0.51</b>	<b>0.74</b>	8.98E7	28.4
	SFP(20%)	✓	93.34	93.25	0.09	8.98E7	28.4
	ASFP(20%)	✓	93.34	93.23	0.11	8.98E7	28.4
	Our1(20%)	✓	93.34	93.17	0.17	8.98E7	28.4
	Our2(20%)	✓	93.34	93.25	0.09	8.98E7	28.4
	SFP(40%)	×	93.66 ± 0.28	92.06 ± 0.68	1.60	5.94E7	52.6
	ASFP(40%)	×	93.66 ± 0.28	92.46 ± 0.43	1.20	5.94E7	52.6
	Our1(40%)	×	93.66 ± 0.28	92.67 ± 0.45	0.99	5.94E7	52.6
	Our2(40%)	×	93.66 ± 0.28	<b>92.92 ± 0.39</b>	<b>0.74</b>	5.94E7	52.6
110	PFEC [19]	×	93.53	92.94	0.61	1.55E8	38.6
	MIL [24]	×	93.63	93.44	0.19	-	34.2
	SFP(20%)	×	94.33 ± 0.40	<b>93.86 ± 0.45</b>	<b>0.47</b>	1.82E8	28.2
	ASFP(20%)	×	94.33 ± 0.40	93.62 ± 0.53	0.71	1.82E8	28.2
	Our1(20%)	×	94.33 ± 0.40	93.61 ± 0.56	0.72	1.82E8	28.2
	Our2(20%)	×	94.33 ± 0.40	<b>93.83 ± 0.58</b>	<b>0.50</b>	1.82E8	28.2
	SFP(40%)	×	94.33 ± 0.40	92.91 ± 0.53	1.42	1.21E8	52.3
	ASFP(40%)	×	94.33 ± 0.40	93.52 ± 0.19	0.81	1.21E8	52.3
	Our1(40%)	×	94.33 ± 0.40	93.66 ± 0.19	0.67	1.21E8	52.3
	Our2(40%)	×	94.33 ± 0.40	<b>93.69 ± 0.24</b>	<b>0.64</b>	1.21E8	52.3

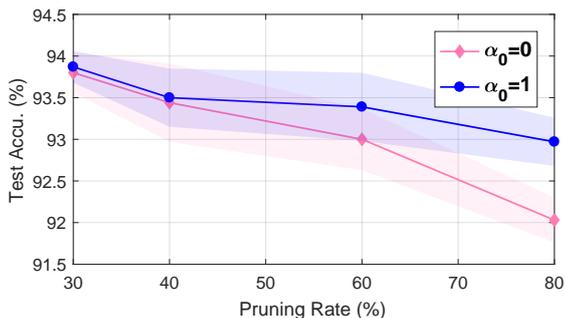


Fig. 3. Transferring to Weight Pruning. Comparison of ResNet-56 on CIFAR-10 between SFP ( $\alpha_0 = 0$ ) and our SRFP ( $\alpha_0 = 1$ ) with the linear decay strategy. (Solid line and shadow represent the mean and standard deviation individually.)

the parameter scheme as [28], and adopt the same data-augmentation scheme as [29].

Our SRFP or ASRFP is adopted after finishing a training epoch. Models are trained from scratch by default. We also provided results with pre-trained models, where the learning rate is one tenth of that of models trained from scratch. The experiments are repeated five times, reported by the “mean  $\pm$  std”. Then we present the results comparing with other state-of-the-art methods, e.g., SFP [8], ASFP [9], MIL [24], PFEC [19], CP [23], ThiNet [27], AutoPruner [22].

### B. ResNet on CIFAR-10

**Settings.** On CIFAR-10, we evaluate our SRFP and ASRFP on ResNet-20/56/110, adopting various pruning rates to study the efficacy of our methods, especially comparing with that

of SFP and ASFP. By default, we use exponential decay with  $\epsilon = 1e - 5$  and  $\alpha_0 = 1$ .

**Results of Filter Pruning.** We conclude the results of both SRFP and ASRFP on CIFAR-10 in Table I. Here, we use “Our1” and “Our2” to refer to SRFP and ASRFP respectively for clarity. We mainly compare our methods with SFP and ASFP. Both SRFP and ASRFP reveal competitive performance on CIFAR-10 compared with other channel pruning techniques across networks of various depths and pruning rates. Notably, our ASRFP performs better than other methods in most cases in Table I. The models are trained from scratch, and the “Accu. Drop” is the accuracy of the pruned model minus that of the baseline model. The smaller is the better.

We compare the results of SFP and SRFP with various pruning rates and network depths across CIFAR-10, shown in Figure 4. When the pruning rate is as small as 20%, the behaviors of the above four methods are quite similar. As the pruning rate increases to 60%, our SRFP and ASRFP outperform SFP and ASFP.

**Transferability to Weight Pruning.** Since filter pruning is a special case of weight pruning, we test the transferability of our SRFP method to weight pruning on CIFAR-10, pruning ResNet-56 with diverse pruning rates, using the linear decay strategy defined by (7). The results are shown in Figure 3, from which we can verify the transferability of our ASFP to weight pruning. Notably, our ASFP has larger relative advantages than SFP in case of larger pruning rates.

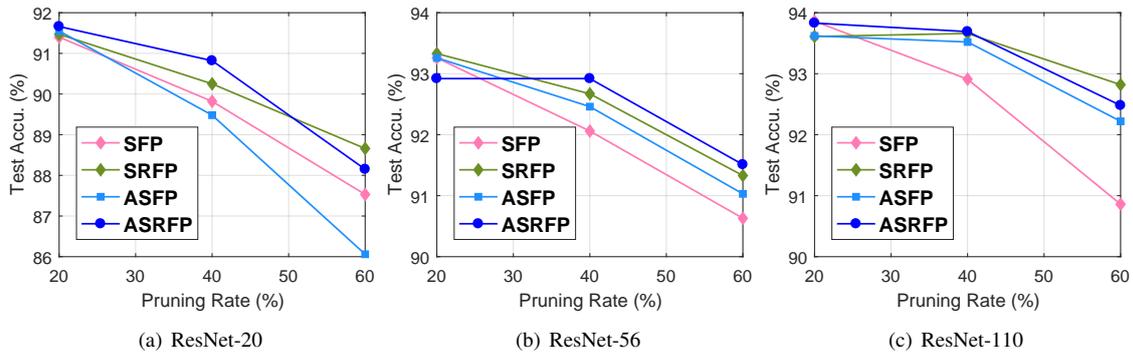


Fig. 4. Comparison of Test Accuracies of ResNet-20/56/110 on CIFAR-10 among SFP/ASFP/SRFP/ASRFP with the pruning rate changing.

TABLE II  
COMPARISON OF PRUNING RESNETS ON IMAGENET

Depth	Method	Pre-trained?	Top-1 Accu.	Top-1 Accu.	Top-5 Accu.	Top-5 Accu.	Top-1 Accu.	Top-5 Accu.	Pruned FLOPs(%)
			Baseline(%)	Accelerated(%)	Baseline(%)	Accelerated(%)	Drop(%)	Drop(%)	
18	MIL [24]	×	69.98	66.33	89.24	86.94	3.65	2.30	34.6
	SFP(30%) [8]	×	70.23	67.10	89.51	87.78	3.13	1.73	41.8
	ASFP(30%) [9]	×	70.23	67.26	89.51	87.88	2.97	1.63	41.8
	Our1(30%)	×	70.23	<b>68.06</b>	89.51	<b>88.06</b>	<b>2.17</b>	<b>1.45</b>	41.8
	Our2(30%)	×	70.23	67.25	89.51	87.59	2.98	1.92	41.8
	ASFP(40%) [9]	×	70.23	65.44	89.51	86.47	4.79	3.04	53.5
	Our2(40%)	×	70.23	<b>65.56</b>	89.51	<b>86.58</b>	<b>4.67</b>	<b>2.93</b>	53.5
34	MIL [24]	×	73.42	72.99	91.36	91.19	0.43	0.17	24.8
	SFP(30%) [8]	×	73.92	71.15	91.62	89.54	2.77	1.67	41.1
	ASFP(30%) [9]	×	73.92	71.17	91.62	90.11	2.75	1.46	41.1
	Our1(30%)	×	73.92	71.35	91.62	90.20	2.57	1.42	41.1
	Our2(30%)	×	73.92	<b>71.39</b>	91.62	<b>90.23</b>	<b>2.53</b>	<b>1.39</b>	41.1
	ASFP(40%) [9]	×	73.92	68.79	91.62	88.95	5.13	2.67	52.7
	Our1(40%)	×	73.92	<b>70.73</b>	91.62	<b>89.87</b>	<b>3.19</b>	<b>1.75</b>	52.7
	Our2(40%)	×	73.92	70.42	91.62	89.63	3.50	1.99	52.7
50	SFP(40%) [8]	×	76.15	73.04	92.87	91.40	3.11	1.47	53.5
	ASFP(40%) [9]	×	76.15	72.98	92.87	91.48	3.17	1.39	53.5
	Our1(40%)	×	76.15	<b>73.62</b>	92.87	<b>91.74</b>	<b>2.53</b>	<b>1.13</b>	53.5
	Our2(40%)	×	76.15	73.60	92.87	91.61	2.55	1.26	53.5
	ThiNet[27]	✓	72.88	72.04	91.14	90.67	0.84	0.47	36.7
	AutoPruner [22]	✓	76.15	74.76	92.87	92.15	1.39	0.72	51.2
	SFP(30%) [8]	✓	76.15	62.14	92.87	84.60	14.01	8.27	41.8
	ASFP(30%) [9]	✓	76.15	75.53	92.87	92.73	0.62	0.14	41.8
	Our1(30%)	✓	76.15	75.98	92.87	92.81	0.17	0.06	41.8
	Our2(30%)	✓	76.15	<b>76.00</b>	92.87	<b>92.90</b>	<b>0.15</b>	<b>-0.03</b>	41.8

TABLE III  
COMPARISON OF THE LINEAR DECAY STRATEGY AND EXPONENTIAL DECAY OF SRFP ON CIFAR-10

Model	Pruned percent(%)	SFP Accu.(%)	Linear Accu.(%)	Exp. Accu.(%)
ResNet-56	30	93.05 ± 0.17	93.03 ± 0.37	<b>93.22</b> ± 0.38
ResNet-56	40	92.06 ± 0.63	<b>92.77</b> ± 0.16	92.67 ± 0.41
ResNet-110	20	93.86 ± 0.45	<b>93.97</b> ± 0.53	93.61 ± 0.56

### C. ResNet on ILSVRC-2012

**Settings.** On ILSVRC-2012, owing to the excellent performance of our ASRFP method on CIFAR-10, we mainly

focus on the evaluation of our ASRFP on ResNet-18/34/50, especially comparing with that of ASFP. By default, we use exponential decay with  $\epsilon = 1e-7$  and  $\alpha_0 = 1$ . For pre-trained models, we let  $\epsilon = 1e-9$ .

**Results.** We summarize the results of both SRFP and ASRFP in Table II. Here, we use "Our1" and "Our2" to refer to SRFP and ASRFP individually. According to Table II, our ASRFP still outperforms other pruning methods in most cases, especially for networks with large pruning rates like 40%. It is worth noting that the relative advantage of our ASRFP method is larger in Top-1 accuracy than in Top-5 accuracy.

**Convergence Analysis.** To illustrate the intrinsic mechanism of SRFP and ASRFP, we compare different Test Accuracy Drops of ResNet-34 on ILSVRC-2012 among

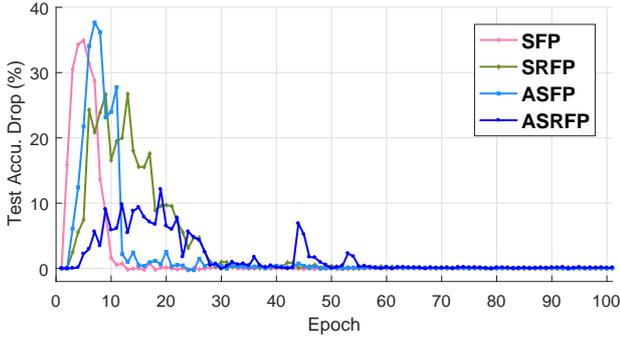


Fig. 5. Comparison of different Test Accuracy Drops of ResNet-34 on ILSVRC-2012 among SFP/ASFP/SRFP/ASRFP with the training epochs increasing when the pruning rate is 30%. The Test Accuracy Drop is the difference between the Top-1 accuracy before pruning and the Top-1 accuracy after pruning, where 0 means that there is no obvious accuracy drops caused by pruning.

SFP/ASFP/SRFP/ASRFP with the training epochs increasing when the pruning rate is 30%, as shown in Figure 5. The Test Accuracy Drop is the difference between the Top-1 accuracy before pruning and the Top-1 accuracy after pruning, where 0 means that there is no obvious accuracy drops caused by pruning.

SRFP, ASRFP and ASFP are all variants of SFP, pursuing better performance at the cost of slowing down the speed of convergence. We notice that the Test Accuracy Drop of SFP converges to 0 at the fastest speed and that of ASRFP converges to 0 at the slowest speed because ASRFP softens the SFP in both pruning rates and the weight decay of pruned filters, thus requiring more epochs to converge, while ASFP and SRFP soften the SFP in pruning rates and the weight decay of pruned filters respectively.

#### D. Ablation Study

We conducted a series of ablation experiments.

**Varying pruning rates and  $\alpha_0$ .** To further shed light on the performance of the SRFP method, we present the results of the estimation of the accuracy of diverse pruning rates and  $\alpha_0$  for ResNet-56/110 in Figure 6(a) and Figure 6(b). Note that SFP is a special case of SRFP with  $\alpha_0 = 0$ . Evidently,  $\alpha_0 = 1$  is a remarkable choice across different pruning rates and network architectures. Besides, as the pruning rate increases, the relative advantage of SRFP with  $\alpha_0 > 0$  is enlarged compared with SFP.

**Different decay strategies.** We compare the results of linear decay and exponential decay on CIFAR-10, as shown in Table III. Both strategies set  $\alpha_0 = 1$ . SFP is a special case of SRFP when  $\alpha_0 = 0$ . We use exponential decay with  $\epsilon = 1e - 5$  for CIFAR-10 and  $\epsilon = 1e - 7$  for ILSVRC-2012 by default. Although linear decay is simple compared with exponential decay and performs well on CIFAR-10, the linear decay strategy produces very poor results on ILSVRC-2012, due to the limited training and pruning epochs. As presented by Figure 7, linear decay strategy decreases the pruned weights in a much smoother and slower manner compared with

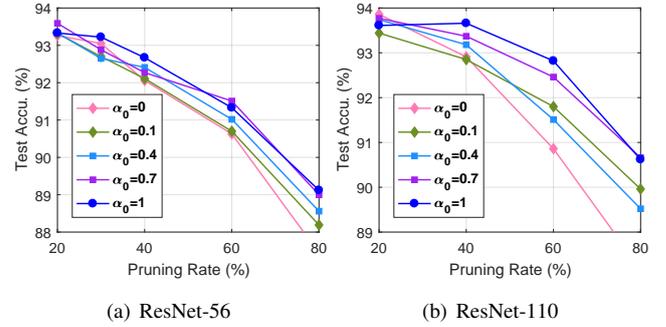


Fig. 6. Comparison of test accuracies of different pruning rates and  $\alpha_0$  for ResNet-56/110 on CIFAR-10 between SFP ( $\alpha_0 = 0$ ) and our SRFP ( $\alpha_0 \neq 0$ ) with the linear decay strategy. SFP is a special case of SRFP with  $\alpha_0 = 0$ .

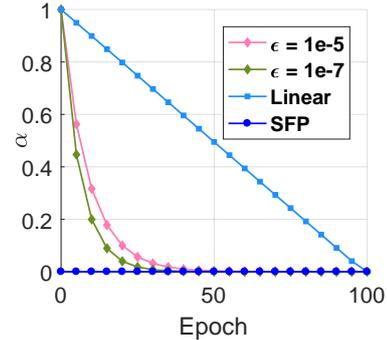


Fig. 7. Comparison of linear decay and exponential decay with various  $\epsilon$  to control the speed of weight decay. All strategies except SFP use  $\alpha_0 = 1$ .

exponential decay, leading to the disastrous non-convergence issue on complex dataset like ILSVRC-2012. And we evaluate different values of  $\epsilon$  for ResNet-18 trained from scratch with the pruning rate of 40% on ILSVRC-2012, including  $1e - 5$ ,  $1e - 7$  and  $1e - 9$ , ultimately choosing  $\epsilon = 1e - 7$ , trading off between the accuracy and convergence.

## V. CONCLUSION

To conclude, we propose a pruning method SRFP and its variant ASRFP, softening the pruning operation of SFP and ASRFP. We present two kinds of weight decay strategies, exponential decay and linear decay and investigate their differences. Our methods perform well across various networks, datasets and pruning rates, also transferable to weight pruning. In theory, our methods are doing the  $\ell_2$ -norm regularization on those pruned filters. Besides, we study the intrinsic mechanism of SRFP and ASRFP and note that SRFP, ASRFP and ASFP pursue better results while slowing down the speed of convergence.

## REFERENCES

- [1] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 2016-December, pp. 770–778, 2016.
- [2] C. Yang, Z. An, H. Zhu, X. Hu, K. Zhang, K. Xu, C. Li, and Y. Xu, "Gated convolutional networks with hybrid connectivity for image classification," in *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020*. AAAI Press, 2020, pp. 12581–12588.

- [3] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2014.
- [4] Y. Zhang, Z. Qiu, T. Yao, D. Liu, and T. Mei, "Fully Convolutional Adaptation Networks for Semantic Segmentation," 2018.
- [5] Z. Liu, J. Li, Z. Shen, G. Huang, S. Yan, and C. Zhang, "Learning Efficient Convolutional Networks through Network Slimming," *Proceedings of the IEEE International Conference on Computer Vision*, vol. 2017-October, pp. 2755–2763, 2017.
- [6] M. Jaderberg, A. Vedaldi, and A. Zisserman, "Speeding up convolutional neural networks with low rank expansions," in *BMVC 2014 - Proceedings of the British Machine Vision Conference 2014*, 2014.
- [7] Z. Liu, M. Sun, T. Zhou, G. Huang, and T. Darrell, "Rethinking the Value of Network Pruning," pp. 1–13, 2018. [Online]. Available: <http://arxiv.org/abs/1810.05270>
- [8] Y. He, G. Kang, X. Dong, Y. Fu, and Y. Yang, "Soft filter pruning for accelerating deep convolutional neural networks," *IJCAI International Joint Conference on Artificial Intelligence*, vol. 2018-July, pp. 2234–2240, 2018.
- [9] Y. He, X. Dong, G. Kang, Y. Fu, C. Yan, and Y. Yang, "Asymptotic Soft Filter Pruning for Deep Convolutional Neural Networks," *IEEE Transactions on Cybernetics*, vol. PP, pp. 1–11, 2019.
- [10] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," *7th International Conference on Learning Representations, ICLR 2019*, 2019.
- [11] J. M. Alvarez and M. Salzmann, "Compression-aware training of deep networks," *Advances in Neural Information Processing Systems*, vol. 2017-December, no. Nips, pp. 857–868, 2017.
- [12] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," *Advances in Neural Information Processing Systems*, vol. 2, pp. 1097–1105, 2012.
- [13] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications," 2017. [Online]. Available: <http://arxiv.org/abs/1704.04861>
- [14] P. Singh, V. K. Verma, P. Rai, and V. P. Namboodiri, "HetConv: Heterogeneous Kernel-Based Convolutions for Deep CNNs," 2019. [Online]. Available: <http://arxiv.org/abs/1903.04120>
- [15] J. Frankle and M. Carbin, "The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks," pp. 1–42, 2018. [Online]. Available: <http://arxiv.org/abs/1803.03635>
- [16] F. Manessi, A. Rozza, S. Bianco, P. Napoletano, and R. Schettini, "Automated Pruning for Deep Neural Network Compression," *Proceedings - International Conference on Pattern Recognition*, vol. 2018-August, pp. 657–664, 2018.
- [17] S. Han, J. Pool, J. Tran, and W. J. Dally, "Learning both weights and connections for efficient neural networks," *Advances in Neural Information Processing Systems*, vol. 2015-Janua, pp. 1135–1143, 2015.
- [18] Y. Guo, A. Yao, and Y. Chen, "Dynamic network surgery for efficient DNNs," *Advances in Neural Information Processing Systems*, pp. 1387–1395, 2016.
- [19] H. Li, A. Kadav, I. Durdanovic, H. Samet, and H. P. Graf, "Pruning Filters for Efficient ConvNets," no. 2016, pp. 1–13, 2016. [Online]. Available: <http://arxiv.org/abs/1608.08710>
- [20] B. O. Ayinde and J. M. Zurada, "Building Efficient ConvNets using Redundant Feature Pruning," pp. 1–9, 2018. [Online]. Available: <http://arxiv.org/abs/1802.07653>
- [21] Z. You, K. Yan, J. Ye, M. Ma, and P. Wang, "Gate Decorator: Global Filter Pruning Method for Accelerating Deep Convolutional Neural Networks," no. NeurIPS, pp. 1–15, 2019. [Online]. Available: <http://arxiv.org/abs/1909.08174>
- [22] J. H. Luo and J. Wu, "AutoPruner: An end-to-end trainable filter pruning method for efficient deep model inference," *Pattern Recognition*, vol. 107, 2020.
- [23] Y. He, X. Zhang, and J. Sun, "Channel Pruning for Accelerating Very Deep Neural Networks," *Proceedings of the IEEE International Conference on Computer Vision*, vol. 2017-October, pp. 1398–1406, 2017.
- [24] X. Dong, J. Huang, Y. Yang, and S. Yan, "More is less: A more complicated network with less inference complexity," in *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, 2017.
- [25] A. Krizhevsky, "Learning Multiple Layers of Features from Tiny Images," ... *Science Department, University of Toronto, Tech. ...*, 2009.
- [26] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge," *International Journal of Computer Vision*, 2015.
- [27] J. H. Luo, J. Wu, and W. Lin, "ThiNet: A Filter Level Pruning Method for Deep Neural Network Compression," *Proceedings of the IEEE International Conference on Computer Vision*, vol. 2017-October, pp. 5068–5076, 2017.
- [28] K. He, X. Zhang, S. Ren, and J. Sun, "Identity mappings in deep residual networks," *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 9908 LNCS, pp. 630–645, 2016.
- [29] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in PyTorch," in *31st Conference on Neural Information Processing Systems (NIPS 2017)*, 2017.