
End-to-End Sensitivity-Based Filter Pruning

Zahra Babaiee¹ Lucas Liebenwein² Ramin Hasani² Daniela Rus² Radu Grosu¹

Abstract

In this paper, we present a novel sensitivity-based filter pruning algorithm (SbF-Pruner) to learn the importance scores of filters of each layer end-to-end. Our method learns the scores from the filter weights, enabling it to account for the correlations between the filters of each layer. Moreover, by training the pruning scores of all layers simultaneously our method can account for layer interdependencies, which is essential to find a performant sparse sub-network. Our proposed method can train and generate a pruned network from scratch in a straightforward, one-stage training process without requiring a pre-trained network. Ultimately, we do not need layer-specific hyper-parameters and pre-defined layer budgets, since SbF-Pruner can implicitly determine the appropriate number of channels in each layer. Our experimental results on different network architectures suggest that SbF-Pruner outperforms advanced pruning methods. Notably, on CIFAR-10, without requiring a pretrained baseline network, we obtain 1.02% and 1.19% accuracy gain on ResNet56 and ResNet110, compared to the baseline reported for state-of-the-art pruning algorithms. This is while SbF-Pruner reduces parameter-count by 52.3% (for ResNet56) and 54% (for ResNet101), which is better than the state-of-the-art pruning algorithms with a high margin of 9.5% and 6.6%.

1. Introduction

Convolutional Neural Networks (CNNs) (LeCun et al., 1989) are widely used in various deep learning computer vision tasks. Large CNNs achieve considerable performance levels, but with significant computing, memory, and energy footprints (Sui et al., 2021). These models are dense and over-parameterized. As a consequence, they cannot be effectively used in resource-limited environments such as mobile or embedded devices. Hence, it's crucial to create smaller models that can perform well without significantly sacrificing their accuracy and performance. This goal can be accomplished either through designing smaller network architectures (Lechner et al., 2020; Tan & Le, 2019) or

End-to-End Sensitivity-based Pruning (SbF-Pruner)

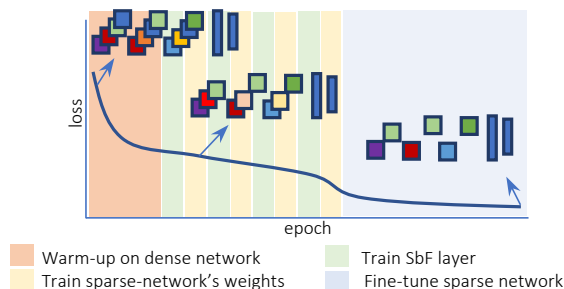


Figure 1. Sensitivity-based filter pruning schedule.

through training an over-parameterized network and sparsifying it by pruning its redundant parameters (Han et al., 2016; Liebenwein et al., 2020; 2021).

Neural network pruning is defined as systematically removing parameters from an existing neural network (Hoeffler et al., 2021). It is a popular technique to reduce the growing energy and performance costs of neural networks and make it feasible to deploy them in resource-constrained environments such as smart devices. Various approaches have been developed to perform pruning as this has gained considerable attention over the past few years (Zhu & Gupta, 2017; Sui et al., 2021; Liebenwein et al., 2021; Peste et al., 2021; Frantar et al., 2021; Deng et al., 2020).

Pruning methods can be categorised into structured and unstructured. Unstructured methods reduce the size by removing individual weight parameters (Han et al., 2016), and structured methods remove parameters in groups, by pruning neurons, filters, or channels (Anwar et al., 2017; Li et al., 2019; He et al., 2018b; Liebenwein et al., 2020). Since our hardware is tuned for dense computations, structured pruning offers a favorable balance between accuracy and performance, providing more computational speedups (Hoeffler et al., 2021).

Filter pruning is a prominent family of structured methods for CNNs. Choosing which filters to remove is the essential part of any filter pruning method. There are several different approaches for sorting the filters by a metric of importance in order to select the least important ones to remove.

Various methods rely on the network structure and weights alone to determine the important filters rather than the training data. Magnitude pruning, for example, is one of the simplest and most common data-free methods, which prunes the filters with the smallest l_1 norm of the weights. Many recent works introduce data-informed methods that focus on the feature maps generated from the training data (or a subset of samples) rather than the filters alone. These methods vary from sensitivity-based approaches, which consider the statistical sensitivity of the output feature maps to the input data (Malik & Naumann, 2020; Liebenwein et al., 2020), to correlation-based methods, with an inter-channel perspective, to keep the least similar or correlated feature maps (Sun et al., 2015; Sui et al., 2021).

Both data-free and data-informed methods often rely on finding the importance of filters of each layer locally. However, the importance of a filter changes relatively to the selection of the filters in previous and next layers. It is challenging to develop methods that can globally determine the importance of filters since CNNs should be considered as a whole. Moreover, a vital element of any pruning method is determining the optimal budget for each layer, which is a problem that all locally-determining importance-metric techniques face. The most trivial way to overcome these issues is to evaluate the network loss with and without each combination of k candidate filters out of N . However, this approach requires the evaluation of $\binom{N}{k}$ subnetworks. It is clear that such computation is not practically feasible.

Training-aware pruning methods aim to learn binary masks to turn on and off each filter. A regularization metric often accompanies them to add a penalty for guiding the masks to the desired budget. Mask learning simultaneously for all filters is an effective method for identifying a globally optimal subset of the network. However, due to the discrete nature of the filters and binary values of the masks, the optimization problem is often non-convex and NP-hard. A simple trick used by many recent works (Gao et al., 2020; 2021; Li et al., 2022) is to use Straight-Through Estimators (Bengio et al., 2013) in order to calculate the derivatives of the binary functions as they are identity functions.

In this paper, we introduce a novel end-to-end pruning technique that simultaneously trains and prunes the network. Instead of learning a binary mask, we train continuous scores from the filter weights with a sensitivity-based filter-pruning (SbF-Pruner) mechanism. This way, we allow the importance scores to be learned via gradient descent to obtain a sparsified network. SbF-Pruner obtains Filter scores through a specialized activation function and automatically computes pruned feature maps as a result of the product of the filter sensitivity scores with feature maps of each layer. Our SbF-Pruner pipeline is symbolically shown in Figure 2.

Our main contributions:

- We introduce the SbF-Pruner, *an end-to-end algorithm that learns the importance scores directly from the network weights and filters*. Our method allows to extract hidden patterns in the filter weights for training the scores, rather than relying only on the weight magnitudes. The feature maps are multiplied by our learned scores during training. This way *our method implicitly accounts for the data samples through loss propagation*, enabling our SbF-Pruner to enjoy the advantages of both data-free and data-informed methods.
- Our SbF-Pruner *automatically calculates global importance scores for all filters and determines layer-specific budgets* with only one global hyper-parameter.
- We empirically investigate the pruning performance of our SbF-Pruner in various pruning tasks and compare it to advanced state-of-the-art pruning algorithms. *Our method proves to be competitive and yield higher pruning ratios while preserving higher accuracy*.

2. End-to-End Sensitivity-Based Filter Pruner

In this section, we first denote our notation and procedurally construct our pruning algorithm.

2.1. Notation

To better describe our approach, we first introduce some useful notation. The filter weights of layer l are represented by $\mathcal{F}_l \in \mathbb{R}^{F \times C \times K \times K}$ where F is the number of filters, C the number of input channels, and K the convolutional kernel size. The feature maps of layer l are given by $A_l \in \mathbb{R}^{F \times H \times W}$ where H and W are the image height and width respectively. For simplicity, we ignore the batch dimension.

2.2. Learning the Scores

The SbF-Pruner can be regarded as an independent layer, whose inputs are the layer filter weights and the outputs are the scores associated with each filter. The pruner layer is a linear transformation of the layer weights to a single vector whose length equals the number of filters of the layer, followed by an activation function which guarantees that the scores are between 0 and $1 + \epsilon$. We will explain in detail the choice of the activation function in the next section.

$$S_l = \phi(\mathcal{F}_l W^{\mathcal{F}}) \quad (1)$$

Here ϕ is the sbf-activation function and the projection is a parameter matrix $W^{\mathcal{F}} \in \mathbb{R}^{(F \times C \times K \times K) \times F}$.

The scores learned by the pruner layer are then pointwise

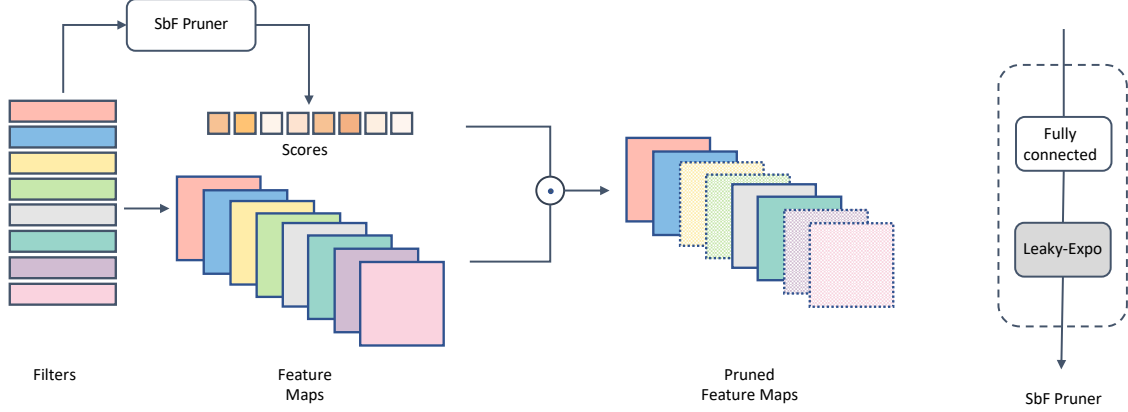


Figure 2. The SbF-Pruner learns the importance scores of the filters from the filter weights.

(\odot) multiplied by the feature maps of the same layer.

$$A'_l = S_l(\mathcal{F}_l) \odot A_l \quad (2)$$

While training the scores, we do not discretize them to binary values; instead, their continuous values are directly multiplied by the feature maps. The closer the filter score is to 1, the more the corresponding feature map is preserved.

2.3. The SbF-Pruner Activation Function

SoftMax is the typical choice of activation function in additive attention for computing importance (Vaswani et al., 2017). However, SoftMax is not a suitable choice for filter scores. While the range of its outputs is between 0 and 1, the sum of its outputs has to be 1, meaning that all scores will have small values, or there would be only one score close to 1.

In contrast to SoftMax, we would intuitively want that many filter scores could possibly be close to 1. More formally, the scores should have the following three main attributes:

1. All filter scores should have a positive value that ranges between 0 and 1, as is the case in SoftMax.
2. All filter scores should adapt from their initial value of 1, as we start with a completely dense model.
3. The filter-scores activation function should have non-zero gradients over their entire input domain.

Sigmoidal activations satisfy Attributes 1 and 3. However, they have difficulties with Attribute 2. For high temperatures, sigmoids behave like steps, and scores quickly converge to 0 or 1. The chance these scores change again is very low, as the gradient is close to zero in these points. Conversely, for low temperatures, scores have a hard time to

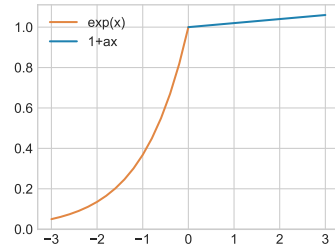


Figure 3. The leaky-exponential function used as the activation function for the SbF-Pruner layers.

converge to 0 or 1. Finding the optimal temperature needs extensive search for each of the layers separately. Moreover, starting from a dense network with all scores set to 1 is not feasible with sigmoids.

To satisfy Attributes 1-3, we designed our own activation function, as shown in Figure 3. First, in order to ensure that scores are positive, we use an exponential activation function, and learn its logarithm values. Second, we allow the activation to be leaky, by not saturating it at 1, as this would result in 0 gradients, and scores getting stuck at 1. Formally, our leaky-exponential activation function is defined as follows, where a is a small value.:

$$\phi(x) = \begin{cases} e^x & \text{if } x < 0 \\ 1.0 + ax & \text{if } x \geq 0 \end{cases} \quad (3)$$

2.4. The Optimization Problem

The network filter-pruning optimization problem can be formulated as in Equation (4), where \mathcal{L} is the loss function, f is the output function of the network with inputs x and labels y , \mathcal{W} are the network parameters, and S are the pruner-layers parameters. By $\|g(S_l)\|_1$ we denote the l_1 -

norm of the binarized filter scores of layer l , and by F_l the number of filters of layer l . Finally, p is the pruning ratio.

$$\begin{aligned} & \min_S \mathcal{L}(y, f(x; \mathcal{W}, S)) \\ \text{s.t.} \quad & \sum_{l=0}^L \|g(S_l)\|_1 - p \sum_{l=0}^L F_l = 0 \end{aligned} \quad (4)$$

To address the pruning budget constraint, we add a regularization factor to the loss function while training the pruner layers. We add the l_1 -norm of the continuous scores calculated for the filters of each layer to the classification loss:

$$\mathcal{L}'(S) = \mathcal{L}(y, f(x; \mathcal{W}, S)) + \lambda \sum_{l=0}^L \|S_l\|_1 \quad (5)$$

Here λ is a global hyper-parameter controlling the regularizer's effect on the total loss. With the loss function consisting of the classification and the l_1 -norm of the scores, the pruning effect of each filter is reflected in the loss function.

By incorporating the classification and l_1 -norm of the scores in the loss function, the effect of the score of each filter is accounted for in the loss value in two ways.

1. When multiplied by small scores, feature maps have less influence on classification. This may result in a larger loss if these maps are important and vice versa.
2. On the other hand, the part of the loss function consisting of the l_1 -norm of the scores, decreases the loss value when there are more scores with small values.

In some sense, this duality mimics the way synapses are pruned or strengthened through habituation or conditioning in nature. If the use of a synapse does not lead to a reward it is pruned, otherwise it is strengthened. Moreover, pruning is essential in order to avoid saturation and enable learning.

The value of the parameter λ plays a very important role in keeping the balance between the two factors, and in controlling the pruning ratio. When λ is set to 0, all scores will stay at values close to 1, and when it is set to a large value, all scores are forced to eventually converge to 0.

During training of the SbF-Pruner layers, we freeze all other network parameters. The scores are directly learned from the filter weights, reflecting attributes such as the filter magnitude and the correlation among the different filters of the same layer. A globally pruned network is obtained by simultaneously training the scores of all layers. This allows to propagate the scored influence of the feature maps of one layer to the feature maps and scores of the next layers.

2.5. The Training schedule

We start training the network by a number of warm-up epochs. During the warm-up phase we train the dense network, with all filter scores fixed to 1 and the pruner layer weights frozen. After the warm-up phase, we start training the pruner layers as well as the network weights.

We multiply the feature maps by the continuous score values. With this approach, when a score is 1, the corresponding feature map is fully preserved. As the score gets smaller, its feature map is getting weaker since it is multiplied by a value smaller than 1, and finally almost completely pruned as the score gets closer to 0. However, while the scores are getting trained, the network weights can learn to adapt to the feature map intensity. This means that, although a filter may have a low score, the loss of the feature map intensity can be compensated by magnifying the feature map itself. In order to prevent this, we use an alternating training approach (Peste et al., 2021).

We train the pruner layers for a few epochs, while all other network parameters are frozen. Then we alternate for training the sparse network weights and during these few epochs we freeze all pruner layer parameters. We multiply the feature maps to their learned scores so far, but we switch to discrete binary values of scores during the sparse network training phases. The binary values are obtained from the continuous scores with a gate function as defined below, where s is the non-binary score value:

$$g(s) = \begin{cases} 1 & \text{if } s \geq 0.5 \\ 0 & \text{otherwise.} \end{cases} \quad (6)$$

After the pruner training phase is completed, we complete the training procedure with fine-tuning the sparse model.

Our proposed method does not necessarily need a pre-trained dense model to find the importance of filters or feature maps. Even when starting from the scratch and after a warm-up phase, it can find an optimal sparse sub-network by training the weights and pruner-layers simultaneously.

2.6. The Algorithm

In previous sections, we introduced the basic concepts of the SbF-Pruner. Algorithm 1 outlines the SbF pruning algorithm for the scores training stage. When training the parameters of the SbF-Pruner layers, Equation (3) is used as the activation function, and during training the rest of the network parameters, we use Equation (6).

2.7. Pruner's Effect on Iterative Inference in ResNets

In this section, we proceed with an analysis of the effect that the SbF-Pruner has on the iterative feature-refinement process in residual networks (Greff et al., 2016).

Algorithm 1 Sensitivity-Based Filter Pruning

Inputs: mini batches \mathcal{B} , Network N parametrized by \mathcal{W} trained for number of warm-up epochs, SBF-Pruner layers parametrized by S , Regularization hyper-parameter λ , Pruning-training epochs E_p , Scores training phase mini-epochs E_s , Fine-tuning epochs E_f .

Output: Pruned Network with Scores for Fine-tuning

```

for  $j$  in  $E_p$  do
  for  $j$  in  $E_s$  do
    for  $b$  in  $\mathcal{B}$  do
      forward calculation:
       $\min_S \mathcal{L}(y, f(x; \mathcal{W}, S)) + \lambda \sum_{l=0}^L \|S_l\|_1$ 
      calculate gradients w.r.t  $S$ .
      Update  $S$  by optimizer.
    end for
  end for
  for  $j$  in  $E_s$  do
    for  $b$  in  $\mathcal{B}$  do
      forward calculation:
       $\min_{\mathcal{W}} \mathcal{L}(y, f(x; \mathcal{W}, S))$ 
      calculate gradients w.r.t  $\mathcal{W}$ .
      Update  $\mathcal{W}$  by optimizer.
    end for
  end for
end for
Return  $N_{\mathcal{W}}$ 
    
```

ResNets (He et al., 2016) are known to enhance representation learning in deeper layers via an iterative feature-refinement scheme (Greff et al., 2016). This scheme suggests that input features to a ResNet do not create new representations. Rather, they gradually and iteratively refine the learned features of the initial residual blocks (Jastrzebski et al., 2018). Iterative refinement of features has been shown to be necessary for obtaining attractive levels of performance with ResNet, while their disruption hurts performance.

As the Sbf-Pruner layer modifies the underlying model structure and the feature maps of a residual block, it is very important to investigate if the Sbf-Pruner preserves the iterative feature-refinement property of ResNets.

To make this analysis precise, let us first formalize iterative inference as discussed in (Jastrzebski et al., 2018): A residual block i in a ResNet with M blocks, performs for the input feature \mathbf{x}_i the following transformation: $\mathbf{x}_{i+1} = \mathbf{x}_i + f_i(\mathbf{x}_i)$. Hence, the following loss function \mathcal{L} can be recursively applied to the network (Jastrzebski et al., 2018):

$$\mathcal{L}(\mathbf{x}_M) = \mathcal{L}(\mathbf{x}_{M-1} + f_{M-1}(\mathbf{x}_{M-1})). \quad (7)$$

A first-order Taylor expansion of this loss, while ensuring that f_j 's magnitude is small, is a good approximation to formally investigating the iterative inference (Jastrzebski

et al., 2018). We thus obtain the following:

$$\mathcal{L}(\mathbf{x}_M) = \mathcal{L}(\mathbf{x}_i) + \sum_{j=i}^{M-1} f_j(\mathbf{x}_j) \cdot \frac{\partial \mathcal{L}(\mathbf{x}_j)}{\partial \mathbf{x}_j} + \mathcal{O}(f_j^2(\mathbf{x}_j)). \quad (8)$$

This approximation reveals that the i -th residual block, modifies features \mathbf{x}_i with roughly the same amount of $f_i(\mathbf{x}_i)$ as that of $\frac{\partial \mathcal{L}(\mathbf{x}_i)}{\partial \mathbf{x}_i}$. This implies a moderate reduction of loss as we transition from the i -th to the M -th block, as an iterative refinement scheme (Greff et al., 2016; Jastrzebski et al., 2018). The refinement step for a vanilla residual block can be computed by the squared norm of $f_i(\mathbf{x}_i)$, and can be normalized to the input feature as: $\|f_i(\mathbf{x}_i)\|_2^2 / \|\mathbf{x}_i\|_2^2$.

Any modification to the structure of the residual network (e.g., if we use the Sbf-Pruner) causes a change in the refinement step. This step has to be investigated if it does or it does not modify the iterative inference scheme.

Lemma 2.1. *The iterative feature-refinement scale is bounded for ResNets with Sbf-Pruner as follows, with parameter ϵ from the leaky integrator and $0 < \delta \leq 1 + \epsilon$:*

$$\frac{\delta}{1 + \epsilon} \frac{\|f_i(\mathbf{x}_i)\|_2^2}{\|\mathbf{x}_i\|_2^2} \leq \frac{\|S_i \odot f_i(\mathbf{x}_i)\|_2^2}{\|S_i \odot \mathbf{x}_i\|_2^2} \leq \frac{1 + \epsilon}{\delta} \frac{\|f_i(\mathbf{x}_i)\|_2^2}{\|\mathbf{x}_i\|_2^2} \quad (9)$$

Proof. A ResNet block i that is equipped with an Sbf-Pruner layer, transforms the features \mathbf{x}_i with the following expression: $\mathbf{x}_{i+1} = S_i \odot (\mathbf{x}_i + f_i(\mathbf{x}_i))$, where S_i stands for the score vector computed by the Sbf-Pruner. The refinement step is given by $\|S_i \odot f_i(\mathbf{x}_i)\|_2^2$ and its input-normalized representation is $\|S_i \odot f_i(\mathbf{x}_i)\|_2^2 / \|S_i \odot \mathbf{x}_i\|_2^2$.

Deriving the upper bound: Assuming that every element in S_i is between δ and $1 + \epsilon$, for $0 < \delta \leq 1 + \epsilon$, we have:

$$\|S_i \odot f_i(\mathbf{x}_i)\| \leq \|S_i\|_{\infty} \|f_i(\mathbf{x}_i)\| \leq (1 + \epsilon) \|f_i(\mathbf{x}_i)\| \quad (10)$$

$$\|S_i \odot \mathbf{x}_i\| \geq \|S_i\|_{\min} \|\mathbf{x}_i\| \geq \delta \|\mathbf{x}_i\|. \quad (11)$$

As a consequence, the following upper-bound inequality holds for the iterative inference:

$$\frac{\|S_i \odot f_i(\mathbf{x}_i)\|_2^2}{\|S_i \odot \mathbf{x}_i\|_2^2} \leq \frac{1 + \epsilon}{\delta} \frac{\|f_i(\mathbf{x}_i)\|_2^2}{\|\mathbf{x}_i\|_2^2} \quad (12)$$

Deriving the lower bound: Assuming that every element in S_i is between δ and $1 + \epsilon$, we have:

$$\|S_i \odot f_i(\mathbf{x}_i)\| \geq \delta \|f_i(\mathbf{x}_i)\| \quad (13)$$

$$\|S_i \odot \mathbf{x}_i\| \leq (1 + \epsilon) \|\mathbf{x}_i\|. \quad (14)$$

As a consequence, the following lower-bound inequality holds for the iterative inference:

$$\frac{\|S_i \odot f_i(\mathbf{x}_i)\|_2^2}{\|S_i \odot \mathbf{x}_i\|_2^2} \geq \frac{\delta}{1 + \epsilon} \frac{\|f_i(\mathbf{x}_i)\|_2^2}{\|\mathbf{x}_i\|_2^2} \quad (15)$$

Inequalities (12) and (15) prove the the stated lemma. \square

Lemma 2.1 has profound implications in practice. It indicates that the iterative-inference property of the ResNets equipped with an SbF-Pruner is both lower and upper bounded. These ResNets not only get compressed in size, but also preserve the representation learning capabilities of ResNets between these two bounds. The bounds themselves can be fine tuned with the parameter λ of Equation (5).

3. Experiments

In this section, we first present our implementation details and then discuss our experimental results.

Baselines. We compare our method to numerous standard and advanced pruning methods. The models include: L1 Norm (Li et al., 2017), Neuron Importance Score Propagation (NISP) (Yu et al., 2017), Soft Filter Pruning (SFP) (He et al., 2018a), Discrimination-aware Channel Pruning (DCP) (Zhuang et al., 2018), DCP-Adapt (Zhuang et al., 2018), Collaborative Channel Pruning (CCP) (Peng et al., 2019), Generative Adversarial Learning (GAL) (Lin et al., 2019), Filter Pruning using High-rank feature map (HRank) (Lin et al., 2020), Discrete model compression (DMC) (Gao et al., 2020), Network pruning via performance maximization (NPPM) (Gao et al., 2021), Channel independence-based pruning (CHIP) (Sui et al., 2021), and Auto-ML for Model Compression (AMC) (He et al., 2018c).

3.1. Implementation Details

Training Procedure. We conducted experiments on CIFAR-10 dataset with two different models, ResNet56 and ResNet110. For both models, we do not use a fully trained network as the baseline to prune. We train each model for 50 warm-up epochs. During the warmup (See Figure 1), we use batch size of 256 and Stochastic Gradient Descent (SGD) as optimizer, with 0.1 as initial learning rate, 0.9 as momentum, and 0.0005 for the weight decay.

After the warm-up stage, we start training all layers of the SbF-Pruner, alternatively, with the rest of network parameters. We train the scores with the SbF-Pruner layers for 3 epochs, while all other parameters are frozen. We add the regularized loss defined in Equation (8) when we train the scores. Then we switch to training the network weights for 6 epochs, while this time the SbF-Pruner layers are frozen and we use the classification loss. In this phase, we use Equation (6) as the gate function for the filter scores. When the SbF-Pruner layers are training, the feature maps get multiplied by the continuous value of the scores. We repeat these two phases for 10 times, summing up to training for 90 epochs. We use the ADAM optimizer with learning rates 10^{-6} and 10^{-3} for training the SbF-Pruners and the network parameters, respectively.

After the scores-training stage is finished, we start fine-

tuning the models. In this stage, we remove all SbF-Pruner layers from the network, keeping the final score vectors trained for each layer. We use Equation (6) to select the binary values of scores, which is our ultimate goal when pruning. Each feature map either is removed entirely (having 0 score), or is completely preserved (having 1 score). We use SGD optimizer with the same parameters as the warm-up stage, and tune the networks for 300 epochs.

Balancing the Pruned Parameters and Flops. When training the SbF-Pruner, we use the regularized loss to guide the scores to our desired budget. However, the number of parameters of each layer is different from the number of flops required for that layer:

$$\begin{aligned} (Params)_l &= C_l \times F_l \times K_l \times K_l \\ (Flops)_l &= C_l \times F_l \times K_l \times K_l \times W_l \times H_l \end{aligned} \quad (16)$$

Since the number of flops is also dependent on the image size, early convolutional layers before the max-pooling layers require more flops as the image sizes are still large. In order to keep the balance between pruned parameters and flops, we multiply the l_1 norm of each layer by the relative input image sizes of that layer to the last layers.

For experiments on ResNet56, we used $\lambda = 5 \times 10^{-4}$ and $\lambda = 15 \times 10^{-4}$ for pruning ratios of 52.3% and 83.0% respectively. On ResNet110, we set λ to 2×10^{-4} and 5.5×10^{-4} to prune 54.9% and 79.3% of the network parameters.

3.2. Performance

Experimental results. Table 3.1 shows our experimental results. We compare the performance of the SbF-Pruner with the state-of-the-art filter pruning methods. As one can see, the SbF-Pruner outperforms other pruning methods, achieving higher accuracy while pruning more parameters. Specifically, on ResNet56 the SbF-Pruner results in an increase in accuracy even compared to the dense baselines while pruning 52.3% of parameters and 49.3% of flops.

The same is true for ResNet110, where the SbF-Pruner achieves an accuracy of 94.69% while it prunes 54.9% of the parameters and 51.3% of the flops. For higher pruning ratios, the SbF-Pruner outperforms CHIP (Sui et al., 2021), the next best method, by 0.37% accuracy increase while pruning 12.8% more parameters, with a total of 83.0% of pruned parameters on ResNet56. Similarly, on ResNet110, the SbF-Pruner achieves higher accuracy than CHIP with 11% more parameters pruned.

Per-layer Budget Discovery. The SbF-Pruner finds the optimal sparse sub-networks in a fully-automated pipeline and does not require budget allocation schedules per layer. Figure 3.2 shows the discovered networks in our experiments from ResNet56 and ResNet110.

In each residual block of the sub-networks learned by the

Table 1. Experiment results on Cifar10 dataset with ResNet56 and ResNet110.

Δ shows the difference in the accuracy of the base dense model used and that of the pruned network. SbF-Pruner does not use a pre-trained base model. Numbers are taken from the reported results in the cited papers.

ResNet56					
Method	Baseline Acc(%)	Pruned Acc(%)	Δ (%)	\downarrow Parameters(%)	\downarrow Flops(%)
l_1 Norm (Li et al., 2017)	93.04	93.06	+0.02	13.7	27.6
NISP (Yu et al., 2017)	N/A	N/A	-0.03	42.2	35.5
SFP (He et al., 2018a)	93.59	93.78	+0.19	N/A	41.1
DCP (Zhuang et al., 2018)	93.80	93.59	-0.31	N/A	50.0
DCP-Adapt (Zhuang et al., 2018)	93.80	93.81	+0.01	N/A	47.0
CCP (Peng et al., 2019)	93.50	93.46	-0.04	N/A	47.0
GAL (Lin et al., 2019)	93.26	93.38	+0.12	11.8	37.6
HRank (Lin et al., 2020)	93.26	93.52	+0.26	16.8	29.3
DMC (Gao et al., 2020)	93.62	93.69	+0.07	N/A	50.0
NPPM (Gao et al., 2021)	93.04	93.40	+0.36	N/A	50.0
CHIP (Sui et al., 2021)	93.26	94.01	+0.75	42.8	47.4
Ours		94.28		52.3	49.3
AMC (He et al., 2018c)	92.80	91.90	-0.90	N/A	50.0
GAL (Lin et al., 2019)	93.26	91.58	-1.68	65.9	60.2
HRank (Lin et al., 2020)	93.26	90.72	-2.54	68.1	74.1
CHIP (Sui et al., 2021)	93.26	92.05	-1.21	71.8	72.3
Ours		92.42		83.0	70.5
ResNet110					
l_1 Norm (Li et al., 2017)	93.53	93.30	-0.23	32.4	38.7
NISP (Yu et al., 2017)	N/A	N/A	-0.18	43.78	43.25
SFP (He et al., 2018a)	93.68	93.86	+0.18	N/A	40.8
GAL (Lin et al., 2019)	93.50	93.59	+0.09	18.7	4.1
HRank (Lin et al., 2020)	93.50	94.23	+0.73	39.4	41.2
CHIP (Sui et al., 2021)	93.50	94.44	+0.94	48.3	52.1
Ours		94.69		54.9	51.3
GAL (Lin et al., 2019)	93.50	92.74	-0.76	44.8	48.5
HRank (Lin et al., 2020)	93.50	92.65	-0.85	68.7	68.6
CHIP (Sui et al., 2021)	93.50	93.63	+0.13	68.3	71.3
Ours		93.68		79.3	74.8

SbF-Pruner, the first layer has lower number of filters remaining after pruning, compared to the second layer. This structure is similar to the bottle-neck architecture present in ResNets with large number of layers. It enables the network to concentrate on the most important features with less capacity, which is exactly what we are looking for with pruning. Many existing pruning algorithms use similar bottle-neck structures, when manually defining layer budgets for pruning (Sui et al., 2021; Lin et al., 2020). The SbF-Pruner is able to discover this pattern automatically, without supervision. Moreover, for high pruning ratios on ResNet110, there are blocks emerging with 0 remaining filters in the first layers. This shows that the SbF-Pruner can also remove entire layers when required, having more freedom in the possible sparse sub-networks search space.

4. Related Work

Filter pruning is a very popular structured method for sparsifying CNNs, which supports storage reduction and processing efficiency without requiring any special library. One can roughly classify existing filter-pruning methods into three main categories based on their selection approaches: Data-free, data-informed, and training-aware methods.

Data-Free Filter Pruning. Following-up on the weights-magnitude-pruning method, where the weights with the smallest absolute values are considered the least important, (Li et al., 2017) use the sum of absolute values of the weights in a filter (the l_1 norm) to prune the filters with the smallest weight values. (He et al., 2018a) dynamically prune the filters with the smallest l_2 norm value in each

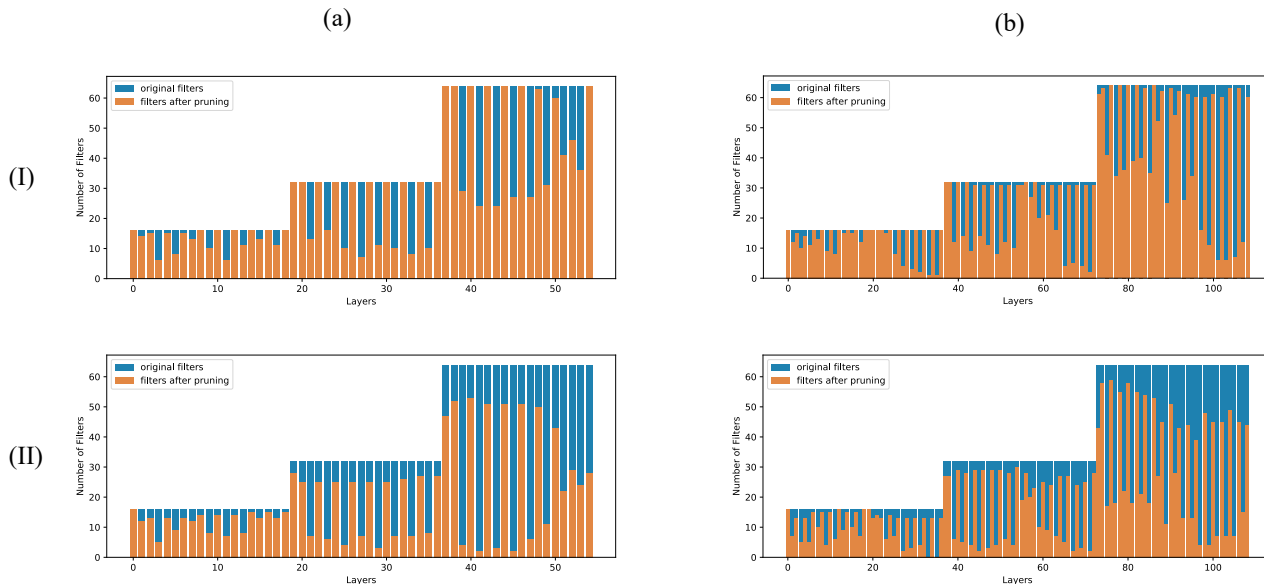


Figure 4. Networks discovered by the SbF-Pruner from a) ResNet56 and b) ResNet110 with medium (I) and (II) high pruning ratios. The SbF-Pruner automatically discovers the optimal per-layer budgets and does not require any pre-defined budgets.

epoch by setting them to zero and repeat this in each epoch during training. (He et al., 2019) use the geometric median of filters as the pruning criterion. Although data-free methods can gain acceptable performance levels, several later works have shown that considering the training data will notably improve pruning precision (Hoefler et al., 2021).

Data-Informed Filter Pruning. Many pruning methods focus on the feature maps since they provide rich information from the data distribution as well as the filters. (Lin et al., 2020) prune the filters whose feature maps have the lowest ranks, and (Liebenwein et al., 2020) use a sensitivity measure to prune filters with lowest effect on the outputs, providing provable sparsity guarantees. Motivated by the importance of inter-channel perspective for pruning, (Sui et al., 2021) use the nuclear norm of the feature maps as an independence metric to prune the filters whose feature maps are the most dependent on the others.

Training-Aware Filter Pruning. These methods use the power of training to learn a filter importance metric or guide the network to a sparse structure. Based on the idea of magnitude pruning, some methods add regularization factors to the loss to directly guide filters to close to zero values. (Wen et al., 2016) and (Louizos et al., 2018) use Group Lasso and l_0 regularization, respectively. Instead of solely relying on the weight magnitudes, (Zhuang et al., 2018) proposes a discrimination-aware channel-pruning method by defining per-layer classification losses. (Gao et al., 2020) train binary gate functions with Straight Through Estimators and (Gao et al., 2021) focus on training binary gates by directly maximizing the accuracy of subnetworks.

No method discussed above however, is able to directly learn the importance-scores from the filter-weights, extract hidden correlations among filters, automatically calculate global importance scores for all filters and determine layer-specific budgets, all at the same time during training, thus taking advantage of both data-free and data-informed methods.

5. Conclusion

We proposed the SbF-Pruner an end-to-end sensitivity-based filter-pruning algorithm that learns importance scores via gradient descent. In contrast to a large spectrum of advanced pruning algorithms, ours does not require a baseline pre-trained network to prune from. It rather sparsifies dense networks from scratch, through cycles of gradient descent.

We showed comprehensively that much better compression rates are achievable through the use of the SbF-Pruner for residual networks, while maintaining a competitively high accuracy. The SbF-Pruner computes global importance scores for filters and automatically associates pruning budgets to a neural network’s layers with a single hyperparameter. This way, the algorithm takes advantage of both the data-free and the data-informed methods.

We hope that future work could begin using our proposed SbF-Pruner framework in resource-hungry applications domains, such as neural architecture search (Mellor et al., 2021), and obtain compressed neural models endowed with salient features, automatically distilled during training, for resource-constraint environments.

References

- Anwar, S., Hwang, K., and Sung, W. Structured pruning of deep convolutional neural networks. *J. Emerg. Technol. Comput. Syst.*, 13(3), feb 2017. ISSN 1550-4832. doi: 10.1145/3005348. URL <https://doi.org/10.1145/3005348>.
- Bengio, Y., Léonard, N., and Courville, A. C. Estimating or propagating gradients through stochastic neurons for conditional computation. *CoRR*, abs/1308.3432, 2013. URL <http://arxiv.org/abs/1308.3432>.
- Deng, L., Li, G., Han, S., Shi, L., and Xie, Y. Model compression and hardware acceleration for neural networks: A comprehensive survey. *Proceedings of the IEEE*, 108(4):485–532, 2020.
- Frantar, E., Kurtic, E., and Alistarh, D. Efficient matrix-free approximations of second-order information, with applications to pruning and optimization. *arXiv preprint arXiv:2107.03356*, 2021.
- Gao, S., Huang, F., Pei, J., and Huang, H. Discrete model compression with resource constraint for deep neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- Gao, S., Huang, F., Cai, W., and Huang, H. Network pruning via performance maximization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9270–9280, June 2021.
- Greff, K., Srivastava, R. K., and Schmidhuber, J. Highway and residual networks learn unrolled iterative estimation. *arXiv preprint arXiv:1612.07771*, 2016.
- Han, S., Mao, H., and Dally, W. J. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding, 2016.
- He, K., Zhang, X., Ren, S., and Sun, J. Identity mappings in deep residual networks. In *European conference on computer vision*, pp. 630–645. Springer, 2016.
- He, Y., Kang, G., Dong, X., Fu, Y., and Yang, Y. Soft filter pruning for accelerating deep convolutional neural networks. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*, pp. 2234–2240. International Joint Conferences on Artificial Intelligence Organization, 7 2018a. doi: 10.24963/ijcai.2018/309. URL <https://doi.org/10.24963/ijcai.2018/309>.
- He, Y., Kang, G., Dong, X., Fu, Y., and Yang, Y. Soft filter pruning for accelerating deep convolutional neural networks, 2018b.
- He, Y., Lin, J., Liu, Z., Wang, H., Li, L.-J., and Han, S. Amc: Automl for model compression and acceleration on mobile devices. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018c.
- He, Y., Liu, P., Wang, Z., Hu, Z., and Yang, Y. Filter pruning via geometric median for deep convolutional neural networks acceleration. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- Hoefler, T., Alistarh, D., Ben-Nun, T., Dryden, N., and Peste, A. Sparsity in deep learning: Pruning and growth for efficient inference and training in neural networks, 2021.
- Jastrzebski, S., Arpit, D., Ballas, N., Verma, V., Che, T., and Bengio, Y. Residual connections encourage iterative inference. In *International Conference on Learning Representations (ICLR)*, 2018.
- Lechner, M., Hasani, R., Amini, A., Henzinger, T. A., Rus, D., and Grosu, R. Neural circuit policies enabling auditable autonomy. *Nature Machine Intelligence*, 2(10): 642–652, 2020.
- LeCun, Y., Boser, B., Denker, J., Henderson, D., Howard, R., Hubbard, W., and Jackel, L. Backpropagation applied to handwritten zip code recognition. *Neural Computation*, 1(4):541–551, 1989.
- Li, H., Kadav, A., Durdanovic, I., Samet, H., and Graf, H. P. Pruning filters for efficient convnets, 2017.
- Li, Y., Gu, S., Gool, L. V., and Timofte, R. Learning filter basis for convolutional neural network compression, 2019.
- Li, Y., Pinteá, S. L., and van Gemert, J. C. Equal bits: Enforcing equally distributed binary network weights. *Proceedings of the AAAI Conference on Artificial Intelligence*, abs/2112.03406, 2022. URL <https://arxiv.org/abs/2112.03406>.
- Liebenwein, L., Baykal, C., Lang, H., Feldman, D., and Rus, D. Provable filter pruning for efficient neural networks, 2020.
- Liebenwein, L., Hasani, R., Amini, A., and Rus, D. Sparse flows: Pruning continuous-depth models. *Advances in Neural Information Processing Systems*, 34, 2021.
- Lin, M., Ji, R., Wang, Y., Zhang, Y., Zhang, B., Tian, Y., and Shao, L. Hrank: Filter pruning using high-rank feature map. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

- Lin, S., Ji, R., Yan, C., Zhang, B., Cao, L., Ye, Q., Huang, F., and Doermann, D. Towards optimal structured cnn pruning via generative adversarial learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- Louizos, C., Welling, M., and Kingma, D. P. Learning sparse neural networks through l0 regularization. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=H1Y8hhg0b>.
- Malik, S. A. and Naumann, U. Interval Adjoint Significance Analysis for Neural Networks. In *Computational Science – ICCS 2020 : 20th International Conference, Amsterdam, The Netherlands, June 3–5, 2020, Proceedings, Part III / edited by Valeria V. Krzhizhanovskaya, Gábor Závadoszky, Michael H. Lees, Jack J. Dongarra, Peter M. A. Sloot, Sérgio Brissos, João Teixeira*, volume 12139 of *Lecture notes in computer science*, pp. 365–378, Cham, Jun 2020. 20th International Conference on Computational Science, Amsterdam (Netherlands), 3 Jun 2020 - 5 Jun 2020, Springer International Publishing. doi: 10.1007/978-3-030-50420-5_27. URL <https://publications.rwth-aachen.de/record/816737>. Weitere Reihe: Theoretical Computer Science and General Issues ; 12139. - Springer eBook Collection.
- Mellor, J., Turner, J., Storkey, A., and Crowley, E. J. Neural architecture search without training. In *International Conference on Machine Learning*, pp. 7588–7598. PMLR, 2021.
- Peng, H., Wu, J., Chen, S., and Huang, J. Collaborative channel pruning for deep networks. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 5113–5122. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/peng19c.html>.
- Peste, A., Iofinova, E., Vladu, A., and Alistarh, D. Ac/dc: Alternating compressed/decompressed training of deep neural networks. *Advances in Neural Information Processing Systems*, 34, 2021.
- Sui, Y., Yin, M., Xie, Y., Phan, H., Zonouz, S., and Yuan, B. Chip: Channel independence-based pruning for compact neural networks, 2021.
- Sun, Y., Wang, X., and Tang, X. Sparsifying neural network connections for face recognition, 2015.
- Tan, M. and Le, Q. EfficientNet: Rethinking model scaling for convolutional neural networks. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 6105–6114. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/tan19a.html>.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. In *Advances in neural information processing systems*, pp. 5998–6008, 2017.
- Wen, W., Wu, C., Wang, Y., Chen, Y., and Li, H. Learning structured sparsity in deep neural networks. In Lee, D., Sugiyama, M., Luxburg, U., Guyon, I., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL <https://proceedings.neurips.cc/paper/2016/file/41bfd20a38bb1b0bec75acf0845530a7-Paper.pdf>.
- Yu, R., Li, A., Chen, C., Lai, J., Morariu, V. I., Han, X., Gao, M., Lin, C., and Davis, L. S. NISP: pruning networks using neuron importance score propagation. *CoRR*, abs/1711.05908, 2017. URL <http://arxiv.org/abs/1711.05908>.
- Zhu, M. and Gupta, S. To prune, or not to prune: exploring the efficacy of pruning for model compression, 2017.
- Zhuang, Z., Tan, M., Zhuang, B., Liu, J., Guo, Y., Wu, Q., Huang, J., and Zhu, J. Discrimination-aware channel pruning for deep neural networks. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *NIPS Proceedings*, volume 2018-December of *Advances in Neural Information Processing Systems*, pp. 875–886. Neural Information Processing Systems (NIPS), 2018. Advances in Neural Information Processing Systems 2018, NIPS 2018 ; Conference date: 02-12-2018 Through 08-12-2018.