

# Effective Model Compression via Stage-wise Pruning

Mingyang Zhang, Xinyi Yu, Jingtao Rong, Linlin Ou

College of Information Engineering, Zhejiang University of Technology

Hang Zhou, People’s Republic of China

linlinou@zjut.edu.cn

## Abstract

*Automated Machine Learning(Auto-ML) pruning methods aim at searching a pruning strategy automatically to reduce the computational complexity of deep Convolutional Neural Networks(deep CNNs). However, some previous work found that the results of many Auto-ML pruning methods cannot even surpass the results of the uniformly pruning method. In this paper, the ineffectiveness of Auto-ML pruning which is caused by unfull and unfair training of the supernet is shown. A deep supernet suffers from unfull training because it contains too many candidates. To overcome the unfull training, a stage-wise pruning(SWP) method is proposed, which splits a deep supernet into several stage-wise supernets to reduce the candidate number and utilize inplace distillation to supervise the stage training. Besides, A wide supernet is hit by unfair training since the sampling probability of each channel is unequal. Therefore, the fullnet and the tinynet are sampled in each training iteration to ensure each channel can be overtrained. Remarkably, the proxy performance of the subnets trained with SWP is closer to the actual performance than that of most of the previous Auto-ML pruning work. Experiments show that SWP achieves the state-of-the-art on both CIFAR-10 and ImageNet under the mobile setting.*

## 1. Introduction

Deep convolutional neural networks(deep CNNs)[16, 34, 10, 18] have achieved outstanding results in many computer vision tasks. However, deep CNNs comes with a huge computational cost, which limits application on embedded devices(i.e. mobile phone).

To expand the application scope of deep CNNs, channel pruning methods were proposed. Traditional channel pruning methods always rely on human-design rules[12, 3]. Recently, inspired by the Neural Architecture Search(NAS), some AutoML-based pruning work[11, 27, 42] has been proposed to automatically prune channels without a human-design mode. Considering a network with 10 layers(each

layer contains 32 channels), the candidates of each layer and the whole network could be 32 and  $32^{10}$ , respectively. Thus, AutoML-based pruning methods can be seen as fine-grained NAS because of more candidates than normal NAS[8, 2, 25] in each layer.

For above mentioned AutoML-based methods, some based on reinforcement learning or evolutionary algorithm[19, 11, 3] are quite time-consuming due to iterative retrain for each pruned network. To reduce the computation in pruning, some AutoML-based work[45, 44, 27, 42] shares the weights through training a supernet for all candidate pruned networks which are called subnets. A typical weight-sharing pruning approach contains three steps: training a supernet by iteratively sampling and updating different candidates, searching the best subnet based on the evolutionary or greedy algorithm and training the best subnet from the scratch. However, Chu[5] considered that the weight-sharing method causes unfull training in the first step since each candidate(subnet) has only a small sampling probability in training. Moreover, unfull training leads to an inaccurate evaluation in the second step, which means some candidates perform well on weight-sharing while bad on training from the scratch. The problem of inaccurate evaluation is particularly obvious in AutoML-based pruning because of more candidates contained in the supernet. Moreover, the width of a supernet also has an impact on the effectiveness of evaluation, which is analyzed in **Sec. 3.1**

To solve the above-mentioned problem, a stage-wise training and searching approach is proposed in this paper. Inspired by [22], we consider to divide a deep supernet into several stage-wise supernet(i.e. ResNet50[10] consists of 4 stages) for reducing the depth of the supernet. Since each stage of the supernet can be trained and searched independently, the candidate number in one stage is an exponential reduction compared with the whole supernet. With a small search space in each stage, the sampling probability of the candidates is raised, which means that each supernet can be fully trained. To alleviate the unfair training result caused by the width of the supernet, both the fullnet and the tinynet, which has the largest and smallest width respectively, are

trained in each iteration. Besides, we present a distributed evolutionary algorithm where each stage can be searched independently in terms of an evolutionary algorithm(EA). The constraints(i.e. FLOPs, latency) for each EA are provided by another EA, called EA manager, where the EA manager searches for the best combination of FLOPs in each stage. Due to small and independent stage-wise search space, each EA can be sped up in a parallel way.

However, the stage-wise supernet in each stage cannot be trained without internal ground truth. To solve such the problem, an existing pre-trained neural network was used to generate stage-wise feature maps that are viewed as ground truth in each stage. Nevertheless, It is also time-consuming to obtain a pre-trained neural network as the teacher network. Besides, the structural difference between the teacher network and the student networks has a strong impact on distillation results. Hence, a stage-wise inplace distillation method is put forward for the fullnet(the largest width subnet) to supervise the learning of subnets. It is worth noting that the fullnet is jointly trained with other subnets. Thus, there does not exist an extra cost for obtaining the fullnet.

Our contribution lies in four folds:

- We propose a stage-wise training and searching pipeline for both channel pruning and NAS. By splitting a CNN into several stages, the number of stage-wise candidates is exponential reduction in contrast to the net-wise candidates. Hence, each candidate obtains full training, which is the essence of accurate evaluation for searching.
- To conveniently provide stage-wise ground truth for each stage, a stage-wise inplace distillation method is presented through the joint training of the fullnet and subnets. Thus, the fullnet can easily supervise the learning of the subnets by offering stage-wise feature maps.
- To accelerate the searching process, a distributed evolutionary algorithm is suggested. Each stage can be searched by the EA with constraints given by an EA manager in a parallel way.
- Experiments shows that the proposed method can enhance the ranking effectiveness of searching and achieve the state-of-the-art in several datasets.

## 2. Related Works

**Neural Architecture Search.** The purpose of neural network structure search is to automatically find the optimal network structure with reinforcement learning(RL)[47, 48], evolutionary algorithm(EA)[31], gradient [25, 39, 4] and parameter sharing methods[30, 2, 8]. RL-based and

EA-based methods need to evaluate each sampled network by retraining them on the dataset, which is time-consuming. The gradient-based method can simultaneously train and search the optimal subnet through assigning a learnable weight to each candidate operation. However, the gradient-based approach causes unfair training results because some candidates obtain more learning resources than others[5]. Moreover, since gradient-based approaches need more memory for training, it cannot be applied to the large-scale dataset. Parameter sharing methods can search on the large-scale dataset by only activating one candidate in each training iteration. Nevertheless, parameter sharing methods cause unfull training results[5]. Unfull or unfair training results will cause an inaccurate evaluation of searching, which means that the best-searched architecture is not the optimal one after retraining. To solve such a problem, Li[22] proposed a blockwisely searching method, which can more fully train each sampled subnets.

**Pruning for CNNs.** Pruning some redundant weights is a prevalent method to accelerate the inference of CNNs. According to the different granularity of pruning, it is divided into weight pruning and channel pruning. For weight pruning, the individual weights in the channel are removed based on some rules[9], which causes unstructured sparse filters and cannot be accelerated directly on most hardware. Therefore, much recent work focus on channel pruning. Channel pruning methods [26, 13, 19, 40, 12] can accelerate the inference of CNNs on general-purpose hardware by reducing the number of filters since the remaining filters are structural. Though the above methods achieve remarkable improvement in the practicality of pruning, it still needs human-designed heuristics to guide pruning.

**Auto-ML pruning.** Recently, inspired by NAS work, AutoML pruning methods[11, 27, 42, 45, 44] have attracted growing interest in automatically pruning for deep CNNs. Different from NAS, the candidate choices are consecutive in the channel pruning task. Compared with pruning methods based on the human-craft rule, AutoML pruning methods aim to search for the best configuration without manual tuning. AMC[11] adopted a DDPG agent to sample a pruned network. And the performance of the pruned networks is evaluated by training from the scratch, which is time-consuming and cannot be applied to a large-scale dataset. MetaPruning[27] trained a PruningNet that can predict weights for any pruned networks, while the parameter amount of the PruningNet is several times of the target network, which leads to unfull train. AutoSlim[42] first trained a slimmable network[45] where the weights between different widths are shared through the supernet, and then searched the best subnet in terms of the greedy algorithm. However, the width of the convolutional layer in each subnet must be the same in training. This leads to the problem that the best subnet achieves the highest accuracy with

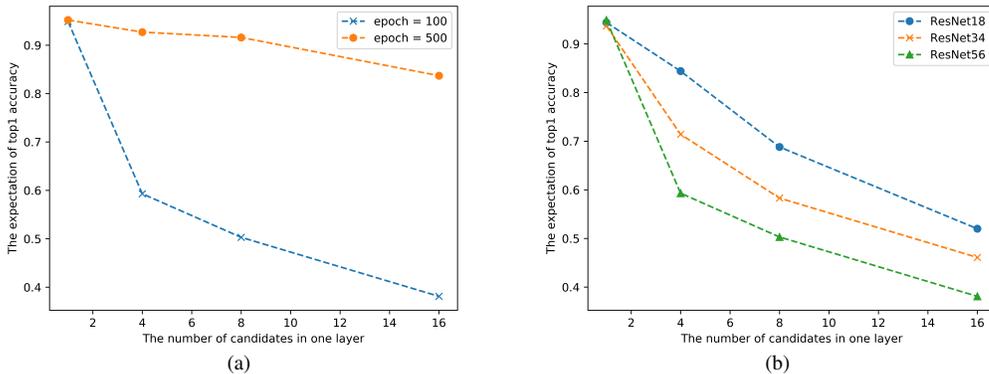


Figure 1. (a) The expectation of top-1 accuracy collected from ResNet56[10] with different candidate number in one layer. The blue and red dash line denotes ResNet56 trained on CIFAR10[35] under 100 epoches and 500 epoches, respectively. (b) The expectation of top-1 accuracy collected from ResNet18[10], ResNet34[10] and ResNet56 with different candidate number in one layer. All models are trained under 100 epoches on CIFAR10[35].

weight sharing, but at the same time, the best subnet gets poor performance when trained from the scratch. To keep the consistency of searching and retraining results, the proposed stage-wise pruning method splits a CNN into several stages and trains them separately under the supervision of the fullnet, which will be explained in Sec.3.

**Knowledge Distillation.** Knowledge distillation is used to train the small student model on a transfer feature set with soft labels or intermediate representations provided by the large teacher model. Soft targets lead to the superior performance of knowledge distillation[14]. However, as the network is designed to become deeper and deeper, it is not enough to just transfer knowledge to a student network from a teacher network by soft targets. To solve such a problem, some previous work[32, 46, 41, 38, 29] transferred the knowledge for the student network from the internal representation of the teacher network. All existing work assumed that the teacher network has been pretrained. Nevertheless, it is always time-consuming to train from scratch to obtain a pretrained teacher network that is essential for various knowledge distillation methods. For example, it may cost more than 10 GPU days to train a ResNet on ImageNet. Moreover, Liu[20] found that the architecture of the teacher and the student networks has a huge impact on transferring results. Hence, we proposed a stage-wise inplace distillation method to overcome the gap and to reduce the time consumption.

### 3. Stage-wise Pruning

The problem of inaccurate evaluation caused by weight sharing is introduced in Sec. 3.1. We find that the depth and width of the supernet have an impact on training adequacy. Thus, the stage-wise inplace-distillation is proposed in Sec. 3.2 to alleviate the aforementioned drawbacks. To effi-

ciently search the optimal subnet from supernet trained by the stage-wise inplace-distillation, we present a distributed evolutionary algorithm in Sec. 3.3.

### 3.1. Challenge of Weight sharing

AutoML pruning methods always need to train a supernet that shares weights for all subnets firstly and almost immediately evaluates the accuracy for each subnet. For many AutoML pruning approaches [11, 27, 42], pruning candidates(subnets) directly compare with each other according to evaluation accuracy. The subnets with higher evaluation accuracy are selected and expected to deliver high accuracy after training from the scratch. However, such an intention cannot be necessarily achieved since some subnets which have outstanding performance on shared parameters perform poorly after training from scratch.

To visualize the performance drop of weight-sharing, we train a supernet with different candidate numbers. For a trained supernet, we randomly sample a batch of subnets from the supernet and evaluate them on the validation dataset. Statistical accuracy expectation  $E(a_{supernet})$  is utilized to evaluate whether the supernet is adequately trained.  $E(a_{supernet})$  is written as

$$E(a_{supernet}) = \sum_{i=1}^n a_{sub_i} \quad (1)$$

where  $n$  denotes the number of the randomly sampled subnet and  $a_{sub_i}$  represents the accuracy of  $i$ th subnet. As shown in Figure 1(a), with the increase of the candidate number, the top-1 accuracy expectation of supernet dramatically degrades under 100 epochs while it falls slightly under 500 epochs. Moreover, we train three different depth supernets and then calculate their  $E(a_{supernet})$ s. It is found that the

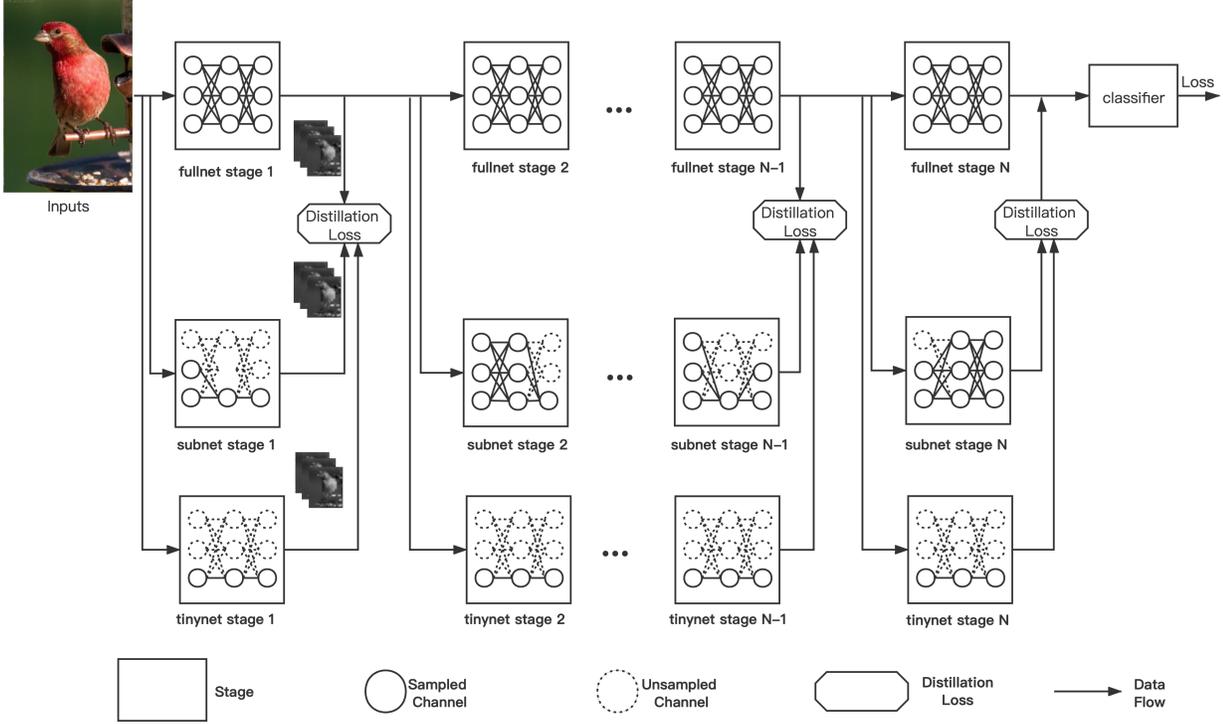


Figure 2. Illustration of the stage-wise training. There are three forms of the network, including the fullnet, the subnet and the tinynet. The fullnet infers inputs once to generate and transfer its knowledge to the subnet and the tinynet by minimizing the L2-distance between the their stage-wise output feature maps. It is worthy to note that these three networks are weight sharing.

expectation of top-1 accuracy is related to the depth of the CNN, which is shown in Figure 1(b). Next, the problem of weight sharing in terms of depth and width is discussed.

**The depth of supernet.** CNNs(i.e. ResNet152[10]) is designed deep to enhance the representative ability, which exponentially increases the subnet number in pruning. The subnet number  $N$  that inherit weights from the supernet can be formulated as

$$\|N\| = g^L \quad (2)$$

where  $g$  denotes the candidate number for each convolutional layer and  $L$  represents the depth of the CNN. For the channel pruning of a deep CNN, the search space  $N$  is always large(e.g.,  $> 30^{50}$ ). Hence, many subnets get unfull training results due to weight sharing, which leads to the ineffectiveness of evaluation.

**The width of layers.** Some Auto-ML pruning work[42, 27, 45, 43, 44] train a single neural network executable at different widths as the supernet. There is not only cross-layer weight sharing but also within the layer. In one layer, the parameters of different widths(candidates) are shared. For instance, all parameters of  $0.25 \times$ (width scaled by 0.25 of the original width) are shared with the half parameters of  $0.5 \times$ . Each sampling on any training step is independent of each other. Thus, for  $n$  training steps, the sampling times

expectation of  $i$ th channel can be formulated as

$$E(c_i) = (1 - \frac{i-1}{m})n \quad (3)$$

where  $m$  denotes the channel number in the layer. Therefore, the channels with small indexes can be trained more times, which causes unfair training results. Formally, we consider a common supernet that contains  $L$  layers, each with  $m$  channels. In pruning, a group of ratio sequence  $R$  can be obtained under certain constraint(i.e. FLOPs), where each ratio sequence  $r = [c_1, \dots, c_L] \in R$ . Because of independent sampling in each layer, the sampling times expectation of  $r$  can be described as

$$E(r) = E(c_1)E(c_2)\dots E(c_L) \quad (4)$$

According to the inequality of arithmetic and geometric means, we have

$$E(r) \leq (\frac{\sum_{i=1}^L E(c_i)}{L})^L \quad (5)$$

Equality holds if and only if  $c_1 = c_2 = \dots = c_L = c$ . Hence,

$$E(r) \leq ((1 - \frac{cm-1}{m})n)^L \quad (6)$$

That is to say, the subnet with uniform sampling obtains the most training resource under certain constraint, which means that the pruning strategy of previous work always prone to get a uniform sampling model.

### 3.2. Stage-wise Inplace Distillation for Training

As mentioned before, too many candidates in training can lead to ineffective evaluation on searching because of unfull training results. To adequately train the supernet, we divide the supernet  $S$  into  $N$  stages according to the depth. Hence, the search space of supernet  $S$  can be represented as

$$S = [S_1, \dots, S_i, S_{i+1}, \dots, S_N] \quad (7)$$

where  $S_i$  denotes the stage-wise supernet of  $i$ th stage. Then we can train the supernet by training the stages separately. The learning of the stage  $i$  can be formulated as

$$W_i^* = \min_{W_i} L_{train}(W_i, S_i; X, Y) \quad (8)$$

where  $X$  and  $Y$  denote the input data and the ground truth labels, respectively. Subsequently, the candidates number for  $i$  th stage can be written as

$$\|S_i\| = g^{L_i} \quad (9)$$

where  $L_i$  denotes the depth of the  $i$ th stage and it is smaller than  $L$ . The search space can be extremely reduced when we train stage-wise supernet in each stage independently. However, internal ground truth in Eq.(8) cannot be obtained directly from the dataset. One solution is using block-wise feature maps generated by a pre-trained network to supervise the training of subnets. However, it is time-consuming to obtain a pre-trained network through training from the scratch in practice(e.g. ResNet50 > 10 GPU days). Besides, the architecture of teacher and student networks has a huge impact on transferring results[20].

To tackle the above problem, the stage-wise inplace distillation is proposed here. The essential idea behind the inplace distillation[44] is to transfer knowledge inside a single supernet from the fullnet to a subnet in each training iteration. For an individual convolutional layer, the performance of the wider candidate cannot be worse than the slim one. Because the wider one can achieve the performance of the slim one by learning weights from some unuseful channels to zeros. Therefore, the performance of any candidates is bounded as follows

$$|y^f - y^s| \leq |y^f - y^r| \leq |y^f - y^s| \quad (10)$$

where  $y^r = \sum_{i=1}^r w_i x_i$  is the aggregated feature,  $r$ ,  $s$  and  $f$  denote the channel number of random sampled, the smallest and the largest candidates, respectively. The rule in Eq.(10) can also be extended to the whole supernet, which means

the performance of the subnet with any width is bounded in the tinynet and the fullnet.

In stage-wise inplace distillation, we use the stage-wise representation of the fullnet to supervise the training of subnets. The pipeline of the supervision with stage-wise inplace distillation is shown in Figure 2. The output  $\hat{Y}_{i-1}$  of the  $(i-1)$ th stage from the fullnet is adopted by the input of the  $i$ th stage of subnets. To supervise the subnets learning from the fullnet, the following MSE loss is considered as the distillation loss in Figure 3

$$L_{train}(\hat{Y}_{i-1}, Y_i) = \frac{1}{F} \|Y_i - \hat{Y}_i\|_2^2 \quad (11)$$

where  $Y_i$  and  $\hat{Y}_i$  denote the output of the subnets and fullnet in  $i$ th stage, respectively,  $F$  is the number of the channels in  $Y$ .

As mentioned in 3.1, the channels with larger indexes suffer from unfull training. To ensure the sufficient training of each channel, an intuitive approach is overtraining. Given a batch of input images and ground truth labels, we first calculate the task loss(e.g. cross-entropy) and the gradients of fullnet through forward and backward propagation, simultaneously. The stage-wise feature maps of fullnet  $\hat{Y} = [\hat{Y}_1, \dots, \hat{Y}_N]$  are saved. Subsequently, under the supervision of  $\hat{Y}$ , the distillation loss in Eq.(11) and the gradients of the stage-wise subnet are calculated. Furthermore, in a subnet training process, we train the smallest width(tinynet) to improve the performance of the supernet. Based on this, each channel can be trained at least once in one iteration. The ground truth label has been generated in the fullnet training process. Thus, the training of each stage-wise subnet can be sped up in a parallel way. The detailed algorithm is described in Algorithm 1. We use **multiprocessing** module<sup>1</sup> to parallelize our algorithm.

### 3.3. Distributed Evolutionary for Searching

After the stage-wise supernet is trained, the learning ability of a subnet can be evaluated by its loss in each stage. However, each stage-wise supernet still contains about  $30^{10}$  stage-wise subnets. It is infeasible to evaluate all of them. For previous one-shot pruning methods, randomly sampling, EA-based and RL-based methods have been used to sample sub-models from the trained supernet for further evaluation. The most recent work found that EA can search better model than RL but need to spend more time on searching. For the proposed stage-wise inplace distillation, a novel method is suggested to search the best subnet according to the stage-wise performance under certain constraint.

Because of inplace distillation training mentioned above,

<sup>1</sup>The multiprocessing module is applied to start the child process and execute our customized tasks in the child process. <https://github.com/python/cpython/tree/3.9/Lib/multiprocessing/>

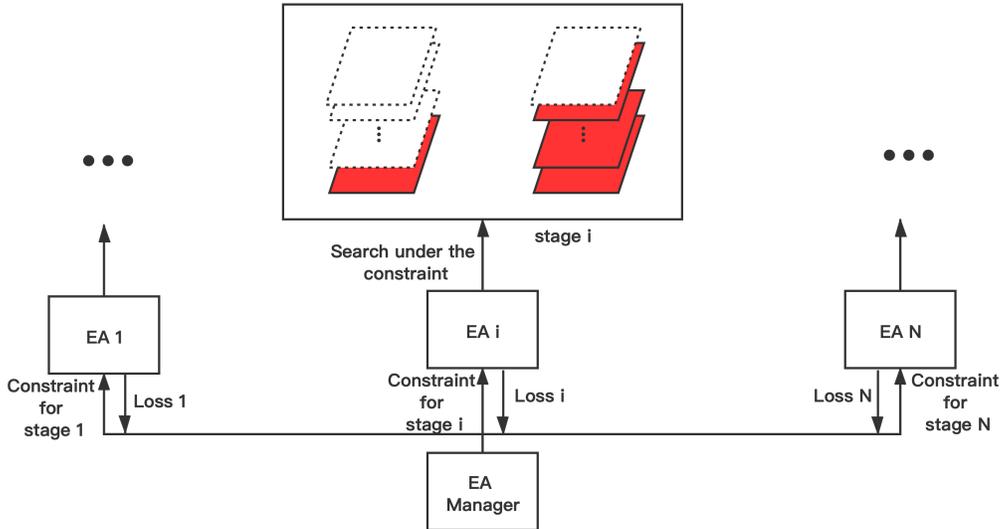


Figure 3. Illustration of the distributed evolutionary algorithm. Both EA Manager and EA are modified from an evolutionary algorithm. Given FLOPs constraint for the whole network, EA Manager is responsible for searching for the best combination of stage-wise FLOPs. The feedback of each FLOPs gene in EA Manager is provided by each EA with searching for the smallest distillation loss under the stage-wise FLOPs constraint.

the EA in Figure 3 is applied to search the best stage-wise subnet that has the smallest distillation loss under given FLOPs constraint. In [27, 48], the genes of each stage-wise subnet were encoded with a vector of channel numbers in each layer. Different from the above work, we aim to imitate the behavior of the teacher in each state. Thus, Eq.(11) is used to evaluate the learning ability of each gene. Because the supernet is split into several stages, the search space of an individual EA is shrunk about  $10^{60}$  times. Then the top  $k$  genes with the lowest loss are selected for generating the new genes with mutation and crossover. The mutation is carried out by changing a proportion of elements in the gene randomly. The crossover means that we randomly recombine the genes in two-parent genes to generate offspring. We can easily enforce the constraint by eliminating the unqualified genes. Through further repeating the top  $k$  selection process and new genes generation process for several iterations, the gene that meets constraints while achieving the lowest loss can be obtained. How can the optimal stage-wise constraints be assigned for each stage?

To automatically find the best assignment plan of stage-wise constraint, a distributed evolutionary algorithm (DEA) is proposed in this section. The workflow of DEA is revealed in Figure 3. The EA Manager is also a kind of evolutionary algorithm that provides the strategy of FLOPs constraint for other EAs. Different from EA above, the genes in EA Manager are encoded with a vector of FLOPs constraint in each stage. The evaluation of each gene is the sum of distillation losses given by all stage-wise EAs. Subsequently,

the top  $k$  genes which are kept generate off-springs in way of mutation and crossover. After several repeating, the optimal stage-wise constraints can be obtained from the top 1 genes. The detail is shown in Algorithm 2. And each stage-wise EA is paralleled to accelerate the searching process. Specifically, we first use the teacher network to generate a batch of representation features for each stage. Therefore, each EA can search for the best stage-wise subnet independently. After searching all of the stages, we can assemble the best model by selecting the best stage-wise subnet from each stage.

## 4. Experiments

In this section, the effectiveness of our proposed stage-wise pruning method is demonstrated. First, We explain the experiment settings on CIFAR-10[35] and ImageNet 2012 dataset[33]. Then, we prune ResNet[10] on CIFAR-10 and visualize the consistency of performance between searching and retraining. Moreover, we apply the stage-wise pruning method to ImageNet 2012 and compare the results with those of other state-of-the-art work. Finally, ablation studies are carried out to find out the influence of using inplace distillation.

### 4.1. Setups

The stage-wise pruning method consists of three steps: **Stage-wise training.** According to resolution size of feature maps, we split ResNet[10] and MobileNet series[16, 34] into 4 and 5 stages, respectively. The distillation loss

---

**Algorithm 1** Framework of supervision with stage-wise in-place distillation.

---

**Require:** The fullnet  $S$ , the stage-wise supernetns  $[S_1, \dots, S_N]$  and the dataset  $(X, Y)$ ;  
**Ensure:** The well-trained stage-wise supernetns  $[S_1, \dots, S_N]$ ;

- 1: **for**  $t = 1, \dots, T$  **do**
- 2:   Get next mini-batch of data  $x$  and label  $y$  from  $(X, Y)$
- 3:   Execute fullnet  $y' = S(x)$ , and save stage-wise feature maps  $\hat{X} = [x, \dots, \hat{Y}_{N-1}]$ ,  $\hat{Y} = [\hat{Y}_1, \dots, \hat{Y}_N]$
- 4:   Calculate cross entropy loss,  $loss = CE(y', y)$
- 5:   Clear gradients,  $optimizer.zero_grad()$
- 6:   Accumulate gradients,  $loss.backward()$
- 7:   Randomly sample width for convolutional layers and obtain stage-wise subnets,  $S_r = [S_{r1}, \dots, S_{rN}]$
- 8:   Uniformly sample smallest width for convolutional layers and obtain stage-wise tinynets,  $S_t = [S_{t1}, \dots, S_{tN}]$
- 9:   **multiprocessing**  $s = [S_r, S_t]$ ,  $x_s = [\hat{X}, \hat{X}]$ ,  $y_s = [\hat{Y}, \hat{Y}]$  **do**
- 10:     Execute subnet,  $y' = s(x_s)$
- 11:     Calculate distillation loss,  $loss = L(y', y_s)$
- 12:     Accumulate gradients,  $loss.backward()$
- 13:   **end multiprocessing**
- 14:   Update weights,  $optimizer.step()$
- 15: **end for**
- 16: **return** trained stage-wise supernetns  $[S_1, \dots, S_N]$ ;

---

of each stage can be calculated by Eq.(11). To match the channel number of the fullnet, the output of each stage is connected with a  $1 \times 1$  convolutional layer without Batch-Norm and non-linear activation. As MetaPruning[27], the width of each convolutional layer is subdivided into 31 ratios from 0.1 to 1.0.

On CIFAR-10[35] dataset, we randomly sample 200 images for each class from training images as validation dataset. The remaining images are used to train the supernet. We use momentum SGD to optimize the weights, with initial learning rate  $\eta = 0.025$ , momentum 0.9, and weight decay  $3 \times 10^{-4}$ . The supernet is trained for 50 epochs with batch size 512 and the learning rate decays  $0.1 \times$  per 10 epochs.

On ImageNet 2012[33] dataset, we randomly sample 50 images for each class from training images as validation dataset. The remaining images are used to train the supernet. We use momentum SGD to optimize the weights, with initial learning rate  $\eta = 0.1$ , momentum 0.9, and weight decay  $3 \times 10^{-4}$ . The supernet is trained for 100 epochs with batch size 512 and the learning rate decays  $0.1 \times$  when the epoch is 30, 60 or 90.

**Stage-wise searching.** After training the stage-wise supernet as above, the best subnet is searched in each stage.

---

**Algorithm 2** Framework of distributed evolutionary algorithm.

---

**Require:** The constraint  $C$ , the fullnet  $S$ , The stage-wise supernetns  $[S_1, \dots, S_N]$  and the dataset  $(X, Y)$ ;  
**Ensure:** The best subnet:  $S_{top}$ ;

- 1: Execute fullnet and save stage-wise feature maps  $\hat{Y}$ ,  $y' = S(X)$ ,  $\hat{Y} = [\hat{Y}_1, \dots, \hat{Y}_N]$
- 2: Randomly generate a batch of genes  $G$  under constraint  $C$ ,  $G = [G_1, \dots, G_N]$ ,  $s.t. ||G_i|| = ||C_{i1}, \dots, C_{iN}|| = C$
- 3: **for**  $t = 1, \dots, T$  **do**
- 4:   **for**  $g = G_1, \dots, G_N$  **do**
- 5:     Obtain stage-wise constraint from  $g$ ,  $g = [C_1, \dots, C_N]$
- 6:     **multiprocessing**  $c = [C_1, \dots, C_N]$ ,  $x_s = X$ ,  $y_s = [\hat{Y}_1, \dots, \hat{Y}_N]$ ,  $s = [S_1, \dots, S_N]$  **do**
- 7:       Search the best stage-wise subnet  $s'$  and calculate distillation loss by EA in Algorithm 3,  $(s', L_{pi}) = EA(s, x_s, y_s, c)$
- 8:     **end multiprocessing**
- 9:     Calculate total loss  $L$  for  $g$ ,  $L = L_{s1} + \dots + L_{sN}$
- 10:   **end for**
- 11:   Keep top k genes  $G_{topk}$  according to  $L$
- 12:   Generate  $M$  mutation genes,  $G_{mutation} = [G_{m1}, \dots, G_{mM}]$ ,  $s.t. ||G_{mi}|| = C$
- 13:   Generate  $H$  crossover genes,  $G_{crossover} = [G_{c1}, \dots, G_{cH}]$ ,  $s.t. ||G_{ci}|| = C$
- 14:   Generate new population  $G$ ,  $G = G_{mutation} + G_{crossover}$
- 15: **end for**
- 16: Select  $S_{top} = [s'_1, \dots, s'_k]$  with smallest  $L$
- 17: **return**  $S_{top}$ ;

---

Firstly, we use the fullnet to generate and save the stage-wise feature maps with 2048 batch size. Subsequently, the hyperparameter of each EA and EA Manager is set to 128 population number, 0.1 mutation probability, 10 iterations. We use 4 and 5 multiprocess to speed up the searching for ResNet and MobileNet series, respectively. Each process can use 2 GPUs to infer with 2048 batch size.

**Retraining** After searching the best subnet, we adopt the same training scheme as [27] on ImageNet 2012 for both ResNet and MobileNet series. The same lines of [25] is followed for the training scheme of ResNet on CIFAR-10. It is noted that all baseline models are trained under the same scheme mentioned above.

## 4.2. Pruning ResNet on CIFAR-10 and Analysis

To demonstrate the effectiveness of stage-wise pruning, we prune ResNet-56[10] under 50% FLOPs constraint on the small dataset of CIFAR-10. As shown in Table 1, our stage-wise pruning method surpasses the baseline model by about 1.4%. Moreover, our method outperforms all other

---

**Algorithm 3** Framework of evolutionary algorithm.

**Require:** The constraint,  $C$ , the stage-wise supernet  $S$ , the stage-wise feature maps  $(X, Y)$

**Ensure:** The best stage-wise subnet:  $S_{top}$  and the stage-wise distillation Loss:  $L$ ;

- 1: Randomly generate a batch of genes  $G$  under constraint  $C$ ,  $G = [G_1, \dots, G_d]$
  - 2: **for**  $t = 1, \dots, T$  **do**
  - 3:   **for**  $g = G_1, \dots, G_d$  **do**
  - 4:     Construct a stage-wise subnet according to  $S$  and  $g$ ,  $S_g$
  - 5:     Calculate the distillation loss of  $S_G$ ,  $L_g = L(S_G(X), Y)$ ,  $L$  from Eq.(11)
  - 6:   **end for**
  - 7:   Keep top  $k$  genes  $G_{topk}$  according to  $L_g$
  - 8:   Generate  $M$  mutation genes under constraint  $C$ ,  $G_{mutation} = [G_{m1}, \dots, G_{mM}]$
  - 9:   Generate  $H$  crossover genes under constraint  $C$ ,  $G_{crossover} = [G_{c1}, \dots, G_{cH}]$
  - 10:   Generate new population  $G$ ,  $G = G_{mutation} + G_{crossover}$
  - 11: **end for**
  - 12: Select  $G_{top}$  with smallest  $L_g$
  - 13: **return**  $G_{top}$ ,  $L_{G_{top}}$ ;
- 

Table 1. Pruning results of ResNet-56.

Method	FLOPs(M)	Top-1 Acc(%)
ResNet-56	125.49	93.27
FP[23]	90.90	93.06
RFP[1]	90.70	93.12
HRank[24]	88.72	93.52
EagleEye[21]	62.23	94.66
<b>SP(Ours)</b>	<b>61.36</b>	<b>95.03</b>

pruning methods in terms of top-1 accuracy.

To evaluate the consistency of model ranking abilities for our method and other AutoML methods, we visualize the relationship between the proxy performance and actual performance. A PruningNet[27] and a Universally Slimming Network(USNet)[44] are trained as supernet under the same training scheme due to fairly compare with MetaPruning [27] and AutoSlim [42]. The total distillation loss is viewed as the proxy performance of our method. The other two methods take the top-1 accuracy of each subnet that inherits weights from supernet as proxy performance. Each subnet will be trained from scratch in order to obtain its actual performance. As shown in Figure 4, the proposed method has a strong correlation between the proxy performance and the actual performance while others barely rank the subnets.

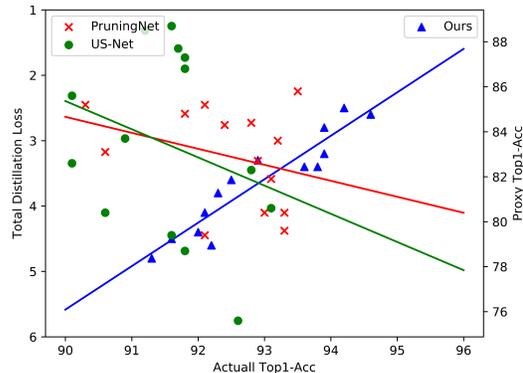


Figure 4. Comparison of ranking effectiveness for Stage-wise Pruning(ours), MetaPruning[27] and AutoSlim[42].

### 4.3. Pruning MobileNet and ResNet on ImageNet 2012

In addition, our method is extended to lighting models on a large-scale dataset, ImageNet 2012. Table 2 summarizes our results on MobileNet V1[16], MobileNet V2[34] and ResNet-50[10]. It is noted that we experiment with both residual and nonresidual networks. We compare our results with uniformly pruned baselines and other recent channel pruning methods. It is shown that our method achieves the best results across different computational budgets. In the case of extreme pruning(i.e. 40M FLOPs), MobileNetV1/MobileNetV2 pruned by SWP outperform baseline model to a considerable degree(9.7% and 6.4%). Figure 5 compares the curve of top-1 accuracy and FLOPs for the most recent Auto-ML pruning methods and uniformly pruning method. Our SWP models can achieve better accuracy with lower computation complexity than other methods.

### 4.4. Pruning under Latency

More and more attention is paid to directly optimize the inference time on the target device. Without knowing the implementation details inside the device, SWP learns to prune channels according to the latency estimated from the device. To evaluate the realistic acceleration, we measure the forward time of the search for the best subnet on one 2080Ti GPU under the latency constraint. The results of MobileNet V1/V2 are shown in 3/4. Under the same compression ratio and similar inference time, our method can obtain better top-1 accuracy.

### 4.5. Combination with NAS

We extend our effective method to NAS. Many NAS work[25, 39, 4] only search operations without considering the channel number. Here, we simultaneously take kernel

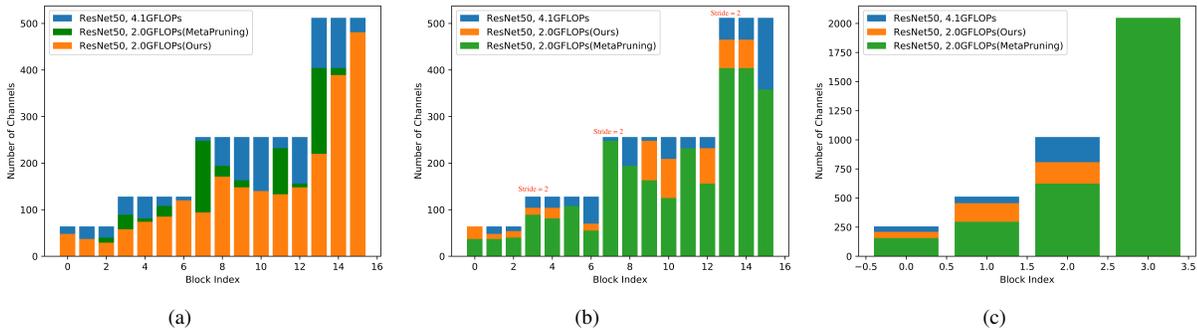


Figure 5. The pruning results of ResNet-50. ResNet-50 is stacked by many blocks which consist of three convolutional layers in the main branch. According to the location, we simply divide the three convolutional layers in each block into top layers, middle layers and bottom layers. (a) The number of channels in top layers. (b) The number of channels in the middle layers. (c) The number of channels in the bottom layers.

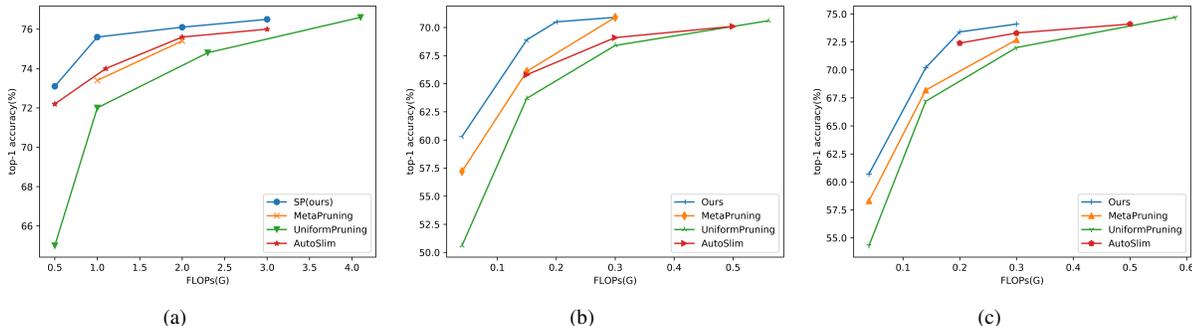


Figure 6. The accuracy-FLOPs tradeoff of Auto-ML pruning methods on (a) ResNet-50, (b) MobileNetV2, (c) MobileNetV1. All models are pruned on ImageNet 2012.

size, layers and channels into account. Our network architectures are designed on MobileNet-V3[15] that consists of a stack with inverted bottleneck residual blocks(MB-Conv). The search space contains the channel number, kernel size and layer number of each stage. The detailed search space is shown in 5.

Several stage-wise supernets are trained on ImageNet dataset[33] using the same settings as MobileNet-V2. After the training, the optimal subnet is found based on distributed evolutionary algorithm. The retraining setting follows MobileNet-V3: RMSProp optimizer with decay 0.9 and momentum 0.9; batch size with 4096; initial learning rate 0.1 that decays by 0.01 every 3 epochs; dropout of 0.8 that only be applied in the fullnet. It is worth noting that a  $3 \times 3$  kernel is central cropped from a  $5 \times 5$  kernel and lower-index layers in each stage are always kept. Therefore, the weights also share on the kernel size and layer size dimension.

The performance of ours and previous SOTA methods with various search spaces are revealed in 6. Our models have 2.0% and 1.7% better top-1 accuracy than MobileNetV3-

small and MobileNetV3-large with similar FLOPs, respectively. Compared with other mobile-setting NAS methods, our method outperforms them under wide-range FLOPs constraints because of more flexible search space. Specifically, our large-size model obtains a competitive result under fewer FLOPs compared with DNA[22] that also split a supernet into several stages but use a pretrained teacher network to supervise. It shows that the architecture of the teacher network has a huge impact on the effectiveness of distillation. We will discuss more details in Sec 4.7.

#### 4.6. Visualization of Searched Models

Channel pruning models are visualized and some insights from the results are discussed. We compare our results with default channels and MetaPruning[27] on ResNet-50. In Figure 6 (a-c), we show the channel number in the top, middle and bottom layers of bottleneck blocks on ResNet-50, respectively. Firstly, it is found that our method is prone to prune more channels from top layers compared with MetaPruning. It is noted that although top layers have a small number of channels, the output

Table 2. Results of ImageNet classification. We show the top-1 accuracy of each method under the same or closed FLOPs.

Network	Method	Acc@1	FLOPs
MobileNet V1	Baseline $0.75\times$	68.4%	325M
	AMC[11]	70.5%	285M
	SN[45]	69.5%	325M
	MP[27]	70.4%	281M
	AutoSlim[42]	69.1%	325M
	<b>SWP(ours)</b>	<b>70.9%</b>	285M
	Baseline $0.25\times$	50.6%	41M
	MP[27]	57.2%	41M
	SN[45]	53.1%	41M
	<b>SWP(ours)</b>	<b>60.3%</b>	41M
MobileNet V2	Baseline $0.75\times$	69.8%	220M
	AMC [11]	70.8%	220M
	MP[27]	71.2%	220M
	SN[45]	68.9%	209M
	AutoSlim[42]	72.4%	207M
	SWP(ours)	<b>73.4%</b>	220M
	Baseline $0.35\times$	54.3%	43M
	MP[27]	58.3%	43M
	<b>SWP(ours)</b>	<b>60.7%</b>	43M
	ResNet-50	Baseline $1.0\times$	76.6%
Baseline $0.75\times$		74.8%	2.3G
SN[45]		74.9%	2.3G
MP[27]		75.4%	2.3G
AutoSlim[42]		75.6	2.3G
AOFP-C1[7]		75.63%	2.58G
C-SGD-50[6]		74.54%	1.7G
ThiNet-50[28]		74.7%	2.1G
<b>SWP(ours)</b>		<b>76.1%</b>	2.0G
Baseline $0.5\times$		72.0%	1.0
SN[45]		72.5%	1.0G
ThiNet-30[28]		72.1%	1.2G
MP[27]		73.4%	1.0G
AutoSlim[42]		74.0%	1.0G
<b>SWP(ours)</b>		<b>75.6%</b>	1.0G

Table 3. Comparison of the realistic inference time with MobileNet V1.

Ratio	Baseline/Pruned time(ms)	Baseline/Pruned Acc@1
$1\times$	0.65/-	70.9%/-
$0.75\times$	0.51/0.51	68.4%/70.2%
$0.5\times$	0.35/0.35	63.3%/67.5%
$0.25\times$	0.21/0.20	49.8%/59.9%

feature maps of the top layer can be extracted by the next middle layer where  $kernel\ size = 3$ . Hence, prune top

Table 4. Comparison of the realistic inference time with MobileNet V2.

Ratio	Baseline/Pruned time(ms)	Baseline/Pruned Acc@1
$1\times$	0.72/-	71.7%/-
$0.75\times$	0.53/0.52	69.8%/71.2%
$0.5\times$	0.41/0.41	65.4%/69.3%
$0.35\times$	0.33/0.32	60.3%/64.8%

Table 5. MobileNetV3-based search space.  $\checkmark$  denotes the search space of the stage containing SE module[17].

Stage	Operator	Ratio	Layers	Kernel Sizes	SE module
	Conv	1.0	1	3	
1	MBCConv	0.1-1.0	1-2	3	$\checkmark$
2	MBCConv	0.1-1.0	2-3	3	$\checkmark$
3	MBCConv	0.1-1.0	2-3	3,5	$\checkmark$
4	MBCConv	0.1-1.0	2-4	3,5	$\checkmark$
5	MBCConv	0.1-1.0	2-6	3,5	$\checkmark$
6	MBCConv	0.1-1.0	2-6	3,5	$\checkmark$
6	MBCConv	0.1-1.0	1-2	3,5	$\checkmark$
	MBCConv	1.0	1	1	

Table 6. Results of NAS methods on Imagenet.

Method	Acc@1	FLOPs
MobileNetV3-large[15]	75.2%	219M
OFA[2]	76.0%	230M
MNasNet[36]	75.2%	315M
DNA-a[22]	77.1%	348M
<b>SWPNas-large(ours)</b>	<b>76.9%</b>	214M
MobileNetV3-small[15]	67.4%	56M
Mnas-small[36]	64.9%	65M
SWPNas-small	69.4%	56M

layers can reduce computational complexity. Secondly, both our method and Metapruning keep more channels for downsampling layers because the feature map size is shrunk. Moreover, our method prunes fewer channels for bottom layers since the feature maps between the subnet and the fullnet should be as close as possible.

## 4.7. Ablation Study

**The choice of the teacher network.** The influence of the distillation strategy in ResNet-50 is analyzed. In our method, the teacher network and the student network are jointly trained by stage-wise inplace distillation. In Strategy 1, we use a pre-trained network that has the same architecture as fullnet to supervise the training of subnets. In

Table 7. Comparison of Stage-wise pruning with different distillation strategies.

Teacher	Student FLOPs	Acc@1
Ours	2.0G	76.1%
	1.0G	75.6%
Strategy 1	2.0G	76.1%
	1.0G	75.6%
Strategy 2	2.0G	75.6%
	1.0G	73.1%

Strategy 2, EfficientNet-B0[37] of which performance surpasses ResNet-50 with lower parameters are employed as the teacher network. The results are shown in Table 7. It is found that the performance of the model searched with inplace distillation method is almost the same as the one searched with a pre-trained method. Hence, the fullnet can supervise the training of the subnets well while training itself. It is unnecessary to spend a lot of time to obtain a pre-trained teacher model. Moreover, despite EfficientNet-B0 has outstanding performance compared with ResNet-50, the models searched with EfficientNet-B0 have worse performance. It may be caused by the large gap between the architectures of the teacher and the student networks.

To further figure out the reason why using ResNet-50 to search can achieve better performance of models than EfficientNet-B0, we visualize the channel number in bottom layers in each block. As shown in Figure 7, the model searched with EfficientNet-B0 keeps fewer channels. And EfficientNet-B0 has much fewer channels than ResNet-50. For example, EfficientNet-B0 has only 40 channels while ResNet-50 has 64 channels in the bottom layer of the first stage. Thus, the student does not need many channels to imitate the stage-wise feature maps generated by EfficientNet-B0. However, over-pruning the channels from the bottom layers will result in insufficient information transmission between the stages.

**The stage number.** In previous experiments, the supernet is simply split into 4(ResNet) or 5(MobileNet series) stages according to downsampling. How many split stages are optimal to transfer knowledge from the teacher to the student? We set the stage number as 8, 4, 2, 1 for ResNet-50. It is seen from Table 8, the performance of the model searched with 8 stages is almost the same as the one searched with 4 stages and surpasses the one searched with lower than 4 stages. Besides, since we need to save the intermediate feature maps of the teacher network to supervise the training of the student network, the more memory is required if larger stage number is set.

Table 8. Comparison of Stage-wise pruning with different stage numbers. We evaluate the models with pruned ResNet-50 under 2.0G FLOPs constraint. The memory is tested under a  $224 \times 224$  image.

Number of stages	Acc@1	Memory
8	76.18%	3M
4	76.12%	1.5M
2	75.88%	0.5M
1	74.53%	0.1M

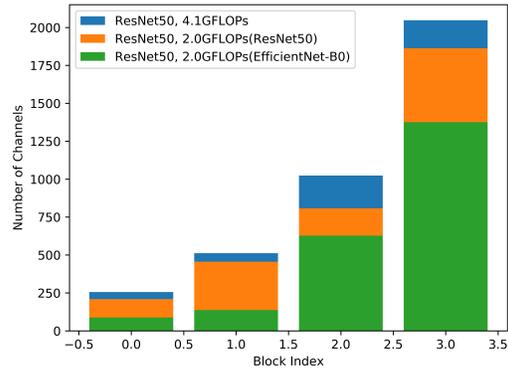


Figure 7. The number of channels in bottom layers with different teachers.

## 4.8. Conclusion

In this work, we have presented stage-wise channel pruning. A stage-wise training process based on inplace distillation and a distributed evolutionary algorithm is proposed in this paper. It is found that the large search space causes low accuracy of the supernet. Hence, we split a supernet into several stage-wise supernets to degrade the complexity of the search space both in training and searching. The experiments show the effectiveness of our proposed method by delivering higher accuracy than the previous work on both CIFAR-10 and ImageNet dataset. The consistency of the proxy and actual subnet performance is greatly improved. Moreover, we extend our method to NAS and also achieve SOTA on the mobile setting. Experiments with various distillation strategies prove that inplace distillation can replace pre-trained distillation, thereby reducing the time to train a teacher network from scratch. We further discuss the impact of stage numbers on search results and found that splitting supernet according to downsampling is the best tradeoff between memory and accuracy.

## Acknowledgement

We gratefully acknowledge the support of National Key R&D Program of China (2018YFB1308400) and

Natural Science Foundation of Zhejiang Province (NO. LY21F030018).

## References

- [1] Babajide O Ayinde and Jacek M Zurada. Building efficient convnets using redundant feature pruning. *arXiv preprint arXiv:1802.07653*, 2018.
- [2] Han Cai, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han. Once-for-all: Train one network and specialize it for efficient deployment. *arXiv preprint arXiv:1908.09791*, 2019.
- [3] Miguel A. Carreira-Perpinan and Yerlan Idelbayev. “learning-compression” algorithms for neural net pruning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [4] Xin Chen, Lingxi Xie, Jun Wu, and Qi Tian. Progressive differentiable architecture search: Bridging the depth gap between search and evaluation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1294–1303, 2019.
- [5] Xiangxiang Chu, Bo Zhang, Ruijun Xu, and Jixiang Li. Fairnas: Rethinking evaluation fairness of weight sharing neural architecture search, 2020.
- [6] Xiaohan Ding, Guiguang Ding, Yuchen Guo, and Jungong Han. Centripetal sgd for pruning very deep convolutional networks with complicated structure, 2019.
- [7] Xiaohan Ding, Guiguang Ding, Yuchen Guo, Jungong Han, and Chenggang Yan. Approximated oracle filter pruning for destructive CNN width optimization. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1607–1616, Long Beach, California, USA, 09–15 Jun 2019. PMLR.
- [8] Zichao Guo, Xiangyu Zhang, Haoyuan Mu, Wen Heng, Zechun Liu, Yichen Wei, and Jian Sun. Single path one-shot neural architecture search with uniform sampling. *arXiv preprint arXiv:1904.00420*, 2019.
- [9] Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. In *Advances in neural information processing systems*, pages 1135–1143, 2015.
- [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [11] 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 *The European Conference on Computer Vision (ECCV)*, September 2018.
- [12] Yang He, Ping Liu, Ziwei Wang, Zhilan Hu, and Yi Yang. Filter pruning via geometric median for deep convolutional neural networks acceleration. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [13] Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [14] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [15] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1314–1324, 2019.
- [16] Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *ArXiv*, abs/1704.04861, 2017.
- [17] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018.
- [18] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [19] Q. Huang, K. Zhou, S. You, and U. Neumann. Learning to prune filters in convolutional neural networks. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 709–718, March 2018.
- [20] Namhoon Lee, Thalayasingam Ajanthan, and Philip H. S. Torr. Snip: Single-shot network pruning based on connection sensitivity. *ArXiv*, abs/1810.02340, 2018.
- [21] Bailin Li, Bowen Wu, Jiang Su, and Guangrun Wang. Eagleeye: Fast sub-net evaluation for efficient neural network pruning. In *European Conference on Computer Vision*, pages 639–654. Springer, 2020.
- [22] Changlin Li, Jiefeng Peng, Liuchun Yuan, Guangrun Wang, Xiaodan Liang, Liang Lin, and Xiaojun Chang. Blockwisely supervised neural architecture search with knowledge distillation, 2020.
- [23] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. *arXiv preprint arXiv:1608.08710*, 2016.
- [24] 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.
- [25] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055*, 2018.
- [26] Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang. Learning efficient convolutional networks through network slimming. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [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. *ArXiv*, abs/1903.10258, 2019.

- [28] Jian-Hao Luo, Jianxin Wu, and Weiyao Lin. Thinet: A filter level pruning method for deep neural network compression, 2017.
- [29] Nikolaos Passalis and Anastasios Tefas. Learning deep representations with probabilistic knowledge transfer. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 268–284, 2018.
- [30] Hieu Pham, Melody Y Guan, Barret Zoph, Quoc V Le, and Jeff Dean. Efficient neural architecture search via parameter sharing. *arXiv preprint arXiv:1802.03268*, 2018.
- [31] Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V Le. Regularized evolution for image classifier architecture search. In *Proceedings of the aaai conference on artificial intelligence*, volume 33, pages 4780–4789, 2019.
- [32] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014.
- [33] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. BergLi, and Fei-Fei. Imagenet large scale visual recognition challenge. In *International Journal of Computer Vision(IJCV)*, December 2015.
- [34] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [35] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.
- [36] Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, Mark Sandler, Andrew Howard, and Quoc V Le. Mnasnet: Platform-aware neural architecture search for mobile. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2820–2828, 2019.
- [37] Mingxing Tan and Quoc V Le. Efficientnet: Rethinking model scaling for convolutional neural networks. *arXiv preprint arXiv:1905.11946*, 2019.
- [38] Hui Wang, Hanbin Zhao, Xi Li, and Xu Tan. Progressive blockwise knowledge distillation for neural network acceleration. In *IJCAI*, pages 2769–2775, 2018.
- [39] Yuhui Xu, Lingxi Xie, Xiaopeng Zhang, Xin Chen, Guo-Jun Qi, Qi Tian, and Hongkai Xiong. Pc-darts: Partial channel connections for memory-efficient differentiable architecture search. *arXiv preprint arXiv:1907.05737*, 2019.
- [40] Jianbo Ye, Xin Lu, Zhe L. Lin, and James Zijun Wang. Rethinking the smaller-norm-less-informative assumption in channel pruning of convolution layers. *ArXiv*, abs/1802.00124, 2018.
- [41] Junho Yim, Donggyu Joo, Jihoon Bae, and Junmo Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4133–4141, 2017.
- [42] Jiahui Yu and Thomas Huang. Autoslim: Towards one-shot architecture search for channel numbers, 2019.
- [43] Jiahui Yu and Thomas S. Huang. Network slimming by slimmable networks: Towards one-shot architecture search for channel numbers. *ArXiv*, abs/1903.11728, 2019.
- [44] Jiahui Yu and Thomas S. Huang. Universally slimmable networks and improved training techniques. *ArXiv*, abs/1903.05134, 2019.
- [45] Jiahui Yu, Linjie Yang, Ning Xu, Jianchao Yang, and Thomas S. Huang. Slimmable neural networks. *ArXiv*, abs/1812.08928, 2018.
- [46] Zhi Zhang, Guanghan Ning, and Zhihai He. Knowledge projection for deep neural networks. *arXiv preprint arXiv:1710.09505*, 2017.
- [47] Barret Zoph and Quoc V. Le. Neural architecture search with reinforcement learning. *ArXiv*, abs/1611.01578, 2016.
- [48] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V. Le. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.