

# Network Pruning via Resource Reallocation

Yuenan Hou<sup>1</sup>, Zheng Ma<sup>2</sup>, Chunxiao Liu<sup>2</sup>, Zhe Wang<sup>2</sup>, and Chen Change Loy<sup>3†</sup>

<sup>1</sup>The Chinese University of Hong Kong <sup>2</sup>SenseTime Group Limited <sup>3</sup>S-Lab, Nanyang Technological University  
<sup>1</sup>hy117@ie.cuhk.edu.hk, <sup>2</sup>{mazheng, liuchunxiao, wangzhe}@sensetime.com, <sup>3</sup>ccloy@ntu.edu.sg

## Abstract

Channel pruning is broadly recognized as an effective approach to obtain a small compact model through eliminating unimportant channels from a large cumbersome network. Contemporary methods typically perform iterative pruning procedure from the original over-parameterized model, which is both tedious and expensive especially when the pruning is aggressive. In this paper, we propose a simple yet effective channel pruning technique, termed network Pruning via rEsouce rEalLocation (PEEL), to quickly produce a desired slim model with negligible cost. Specifically, PEEL first constructs a predefined backbone and then conducts resource reallocation on it to shift parameters from less informative layers to more important layers in one round, thus amplifying the positive effect of these informative layers. To demonstrate the effectiveness of PEEL, we perform extensive experiments on ImageNet with ResNet-18, ResNet-50, MobileNetV2, MobileNetV3-small and EfficientNet-B0. Experimental results show that structures uncovered by PEEL exhibit competitive performance with state-of-the-art pruning algorithms under various pruning settings. Our code is available at <https://github.com/cardwing/Codes-for-PEEL>.

## 1. Introduction

Recent years have witnessed the great success of convolution neural networks in many computer vision tasks, e.g., image classification [5], semantic segmentation [19] and object detection [25]. The remarkable performance usually comes at the cost of dramatically increased network complexity, which makes these models prohibitive for resource-limited edge devices and mobile applications. Channel pruning [14, 20] has been acknowledged as an effective algorithm to notably reduce the model’s demand on computational and storage resources via discarding less informative channels of the original model.

†: Corresponding author.

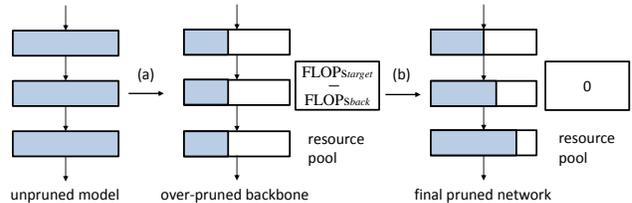


Figure 1: Illustration of PEEL. Given the original unpruned model and target  $FLOPs_{target}$ , PEEL performs the following operations: (a) construct an over-pruned backbone model whose  $FLOPs_{back}$  is smaller than the target value  $FLOPs_{target}$  and store the extra FLOPs ( $FLOPs_{target} - FLOPs_{back}$ ) in a resource pool, (b) reassign the resources in the pool to different layers of this backbone model based on the estimated layer importance and obtain the desired compact model. The blue part in each block denotes the remaining channels in a layer and the white part represents the pruned channels.

Conventional channel pruning approaches [16, 3, 29] typically remove redundant channels from a large, cumbersome neural network to acquire a compact model. For instance, Guo *et al.* [3] introduce additional parameters to learn the retaining probability of each individual channel and prune unnecessary channels based on the learned probability value. Liu *et al.* [16] add  $L_1$  regularization of Batch Normalization (BN) statistics to the training process and recursively eliminate less important channels according to the magnitude of BN statistics. These channel pruning techniques either entail extra learnable parameters or involve iterative and expensive pruning procedure, which prohibits many of them from real-world applications.

In this paper, we propose a conceptually simple yet effective channel pruning approach to obtain a compact and good-performing network with negligible cost. Our algorithm is based on the premise that all layers do not contribute equally to the final performance of the network. Layers that have indispensable effect on the network behavior are regarded as important while other layers are considered as less crucial. By shifting resources from less crucial layers to more significant layers, one can easily acquire a desired

slim model from off-the-shelf network architectures.

Our method is known as network pruning via resource reallocation (PEEL). As the name implies, PEEL addresses channel pruning via reallocating the resources on a well-established backbone model. As shown in Fig. 1, we first construct an over-pruned backbone model and store the saved resources (*e.g.*, FLOPs, parameters or latency) in a resource pool. By estimating layer importance via some standard criteria, resource reallocation is then performed to assign resources to different layers based on the estimated layer significance in one round.

PEEL offers the possibility of reassigning parameters on a backbone architecture. As opposed to previous backward elimination methods, PEEL is built upon forward selection on a predefined architecture. It does not introduce additional learnable parameters and is free from expensive iterative pruning procedure. Despite its simplicity, PEEL can find good structures based on diverse backbone models on ImageNet, and its performance is on par with the architectures pruned by state-of-the-art channel pruning algorithms (*e.g.*, DMCP [3]) under various pruning levels.

In summary, we contribute a simple yet effective channel pruning algorithm, PEEL, to produce a desired compact model via performing resource reallocation on a predefined backbone model. Through extensive experiments, we verify that PEEL can produce stable searching results with small performance variance. Besides, it is applicable to various backbone models and is robust to the criterion to determine the layer importance. We validate the effectiveness of PEEL on modern CNN architectures, *i.e.*, ResNet-18, ResNet-50, MobileNetV2, MobileNetV3-small and EfficientNet-B0, on the large-scale ImageNet benchmark. Under diverse FLOPs constraints, PEEL can consistently find targeted slim models with less searching cost compared with the state-of-the-art channel pruning techniques, *e.g.*, DMCP [3], FPGM [7] and USNet [29]. We hope PEEL can serve as a strong baseline to facilitate future research on channel pruning.

## 2. Related Work

**Channel Pruning:** Channel pruning is a prevailing approach to remove redundant channels from the over-parameterized neural networks. Typically, channel pruning is composed of three stages, *i.e.*, training the original network, pruning unimportant channels and finetuning the pruned architecture [16, 4].

According to the pruning directions, contemporary channel pruning algorithms can be categorized into two types, *i.e.*, backward elimination from the original model [3, 2, 4, 6, 7, 16, 28] and forward selection from an empty network [27, 32]. The majority of the pruning methods fall into the former type and they delete unnecessary connections based on diverse pruning criteria, *e.g.*, the magnitude

of kernel weights [4] or BN statistics [16], the feature map reconstruction errors [20] or the absolute first-order gradient of objective function [21]. For instance, Liu *et al.* [16] add  $L_1$  regularization to the network training phase to enforce network sparsity and remove channels whose BN statistic is smaller than a preset threshold. Guo *et al.* [3] phrase channel pruning as a Markov process and introduce additional parameters to each channel to learn the probability of pruning. Contrary to the preceding indirect pruning algorithms, another line of work [28, 29] concentrates on training a supernet and directly takes the accuracy of the sampled structures as the measurement. For example, Yu *et al.* [28] first train a universally slimmable super-network and then exhaustively evaluate a collection of sampled sub-networks from the trained supernet to determine the optimal substructure. Nonetheless, the aforementioned work suffers from expensive training and searching cost especially when the pruning is aggressive since either iterative pruning procedure or many rounds of evaluation are required to find the best-performing architecture.

As to the forward selection, Ye *et al.* [27] start from an empty network and gradually add important neurons from the original model in the greedy manner. Zhuang *et al.* [32] recursively place important channels in each empty layer based on the gradient norm. Our algorithm belongs to the forward selection category but has several distinctions with [27, 32]. First, we start from a predefined backbone instead of an empty architecture therefore making our selection process more efficient. Second, in contrast to previous efforts that add informative neurons one by one, PEEL first estimates the importance of different layers and then reallocate resources to these layers via the estimated importance in one round, which substantially relieves the selection burden.

Our work is related to network rejuvenation [22]. However, the ultimate objective of [22] is to increase the resource utilization ratio and it is not dedicated to the network pruning task. In all, our PEEL differs from their method in four aspects. First, PEEL reassigns resources according to the computed layer importance while [22] allocates the resources based on the number of remaining channels in each layer. The allocation policy of [22] heavily hinges on the quality of the backbone model and may suffer a lot if the architecture of the backbone model is irrational. Second, [22] requires sophisticated procedure, *e.g.*, parameter reinitialization, neural rescaling and well-designed training scheme, to train the resulting model while PEEL just entails the vanilla training strategy. Third, [22] has several hyperparameters to tune, *e.g.*, the sparsity coefficient, the pruning threshold for each layer and the utilization threshold to conduct neural rejuvenation, while PEEL has only one hyperparameter. Fourth, [22] conducts experiments only on mild pruning levels on ImageNet whereas PEEL includes mild, medium and aggressive pruning settings.

**Knowledge Distillation:** Knowledge distillation [8] is widely acknowledged as an effective technique to transfer the dark knowledge from the large over-parameterized network to the small compact model. The large network is called teacher and the small model is dubbed student. The transferred knowledge can take on many forms, *e.g.*, softened output logits [8], visual attention maps [30, 9] or intermediate feature maps [15, 11]. The above distillation algorithms typically consume huge memory storage since they need to simultaneously load both student and teacher during training. To circumvent the participation of cumbersome teacher models, recent studies [24, 10, 31] have focused on deriving the supervision signals from the student model itself. For instance, Hou *et al.* reinforce the representation learning of shallow layers via attention maps of deep layers [10]. Our algorithm only adopts the simplest form of knowledge distillation [8] to compare fairly with other pruning algorithms. Nevertheless, the aforementioned distillation approaches can also be employed to further boost the performance of the final compact model.

### 3. Methodology

Channel pruning can be formulated as follows: given the target resource constraint (*e.g.*, FLOPs, parameters, latency), the objective is to find a set of channel configurations that achieves the optimal performance on the target task. Considering that the possible number of channels in each layer is a non-negative integer, we phrase channel pruning as a combinatorial optimization problem. The specific mathematical formulation is presented below:

$$\begin{aligned} \max_{(c_1, \dots, c_N)} & f(c_1, \dots, c_N). \\ \text{s.t.} & e(c_1, \dots, c_N) \leq M \\ & 0 \leq c_i \leq C_i \\ & 0 \leq i \leq N \end{aligned} \quad (1)$$

Here,  $f$  is the performance estimation function of the model,  $e$  is the resource estimation function,  $M$  is the target resource constraint (*e.g.*,  $\text{FLOPs}_{target}$  in Fig. 1),  $c_i$  and  $C_i$  are the channel number in the  $i$ -th layer of the pruned and original unpruned model, respectively,  $N$  is the number of layers in the unpruned network.

Conventional approaches [29, 3, 17, 16] typically discard unnecessary channels from the original over-parameterized network. Distinct from previous backward elimination methods, the proposed PEEL addresses channel pruning via reallocating parameters in a predefined architecture, usually an over-pruned backbone from the original network. The resource consumption of this predefined structure is smaller than the target value and the remaining resources are stored in a place named resource pool (Fig. 1). The intuition of PEEL is that those informative layers should be assigned

more parameters to amplify their positive effect on the final performance.

The proposed PEEL is comprised of five steps, namely, the construction of the over-pruned backbone and resource pool, estimation of layer importance, layer grouping, reallocation of parameters to the backbone based on the estimated importance, and finally, the retraining the resulting compact model with knowledge distillation.

**Step 1 - Construction of the over-pruned backbone and resource pool:** Given the target resource constraints  $M$ , we first build an over-pruned backbone model. In this paper, we adopt uniform pruning [18] to construct the backbone model. Uniform pruning possesses many attractive properties. First, it preserves the original architecture as much as possible and therefore the pruned network shares more structural similarity with the original model. Such structural resemblance will facilitate the knowledge distillation process and can bring more gains to the pruned model. Second, contrary to previous pruning methods that require the training of the original unpruned network, uniform pruning introduces no training and searching cost since the pruning ratio is same for all layers. Third, when the pruning is mild, *e.g.*, pruning fewer than 40 % channels, uniform pruning can serve as a strong baseline, with its performance comparable with that of contemporary pruning algorithms. Consequently, we employ uniform pruning to construct a simple backbone model.

We next prepare the resource pool. We define the proportion of the resources the backbone model consumes (*e.g.*,  $\text{FLOPs}_{back}$  in Fig. 1) to the given target resource constraint  $M$  as  $\lambda$  ( $0 \leq \lambda \leq 1$ ), which is a tunable hyperparameter in our algorithm. The resource pool thus stores the extra resources of  $(1 - \lambda)M$ , *e.g.*,  $\text{FLOPs}_{target} - \text{FLOPs}_{back}$  in Fig. 1. Intuitively, a high  $\lambda$  suggests few resources will be kept in the resource pool. With a small  $\lambda$ , more resources will be stored in the pool by having a much slimmer over-pruned backbone. In our experiments,  $\lambda$  is empirically set as 0.8 since this value yields satisfactory results. We avoid an overly small  $\lambda$  observing that over-aggressive pruning will hurt the backbone’s performance, which could hardly be compensated by subsequent operations. It is noteworthy that although we leverage uniform pruning to establish the backbone model, PEEL can still work well when taking different models as backbones.

**Step 2 - Estimation of layer importance:** One crucial drawback of uniform pruning is that it considers all layers to be of equal importance and forces the same pruning ratio for all layers. This practice neglects the fact that some layers contribute more than other layers to the model performance, therefore these influential layers should be assigned more parameters. Hence, one core objective of PEEL is to distinguish important layers from less informative layers and place more parameters in these informative layers.

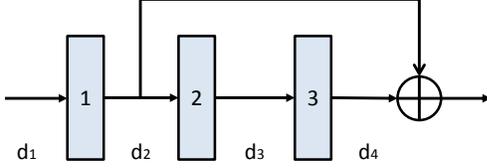


Figure 2: Illustration of pruning channels in residual blocks. For layer  $i$ , the number of input and output channels are  $d_i$  and  $d_{i+1}$ , respectively. The output of layer 1 and 3 are added element-wise via the residual connection. Hence,  $d_2$  and  $d_4$  have to be kept equal after pruning.

There are several options for the evaluating criteria on estimating layer importance, *e.g.*, weight magnitude [4], reconstruction errors of layer activations [20], BN statistics [16], and gradient of the cost function [21]. We eventually choose the BN statistics, *i.e.*, absolute gamma value, as the criterion to reflect the importance of each layer to the final performance as it yields appealing results and requires little computation cost. Batch normalization [13] has been a basic block in modern CNN architectures to address the internal covariate shift. Suppose  $X_i$  and  $X_o$  are the input and output of the BN layer, respectively. BN layer performs the following transformation on the input:

$$X_o = \gamma \frac{X_i - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta, \quad (2)$$

where  $\mu$  and  $\sigma$  are the mean and standard deviation of the current input along the channel dimension, respectively,  $\epsilon$  is a small value to ensure numerical stability,  $\gamma$  and  $\beta$  are two learnable parameters for the BN layer. Following [16], we add the  $L_1$ -norm of  $\gamma$  to the main loss to encourage sparsity:

$$\mathcal{L}_{train} = \mathcal{L}_{cls}(\mathbf{O}_S, \mathbf{O}_L) + \alpha_s \sum_{i=1}^{N_\gamma} |\gamma_i|, \quad (3)$$

where  $N_\gamma$  is the number of all  $\gamma$  in the network,  $\alpha_s$  is the sparsity coefficient,  $\mathbf{O}_S$  is the model prediction,  $\mathbf{O}_L$  is the ground-truth label and  $\mathcal{L}_{cls}$  is the cross entropy loss.

**Step 3 - Layer grouping:** Considering the fact that modern CNN architectures, *e.g.*, ResNet, have residual connections, pruning channels in these networks is troublesome since the output of one layer is directly added to that of another layer, and the channel dimension of these layers has to be kept the same. An example is presented in Fig. 2. To prune these three layers, we have to keep  $d_2$  equal to  $d_4$ . Moreover, ResNet typically has many residual blocks, thus posing more constraints on the consistency in the channel dimensions of multiple layers. To address the aforementioned issue, we decide to group convolutional layers according to their positions in the original network. More concretely, layers whose feature maps have the same spatial size are grouped together. Under this condition, the output channels of the layers in one group will undergo the same

linear expansion or shrinkage, which satisfies the dimension consistency requirement raised by residual connections.

**Step 4 - Resource reallocation:** After the layer grouping process, we need to properly assign the remaining resources in the resource pool to those groups of layers. There are two options for the resource reallocation procedure, *i.e.*, allocating the additional parameter resources in one round or in multiple rounds. The primary distinction between these two options is that the former only considers layer importance in the current uniformly pruned model while the latter takes into account the ever-changing layer significance in various intermediate models. Apparently, the one-round practice will consume much less computation budget than the multiple-round version. Here, we focus on the one-round allocation as it can already bring appealing results and leave the comparison between one-round and multi-round assignment in the ablation study.

Suppose there are  $G$  groups in total after layer grouping and the total number of their output channels in group  $i$  is  $g_i$ . Then, for group  $i$ , its importance is calculated as:

$$T_i = \frac{1}{g_i} \sum_{j=1}^{g_i} |\gamma_j|. \quad (4)$$

Then, the resources assigned to group  $i$  are:

$$R_i = \frac{T_i}{\sum_{j=1}^G T_j} (1 - \lambda)M, \quad (5)$$

where  $(1 - \lambda)M$  is the resources available in the resource pool. In each group, we adopt the linear expansion strategy, *i.e.*, the channels of all layers in one group are multiplied by a common value. In other words, the increased channels of one layer is proportional to its original number of channels. After performing the linear expansion on each layer, we eventually obtain the desired compact model.

**Step 5 - Distilling knowledge from the original model:**

Once the compact model is obtained after resource allocation, we follow the common practice in network pruning [3, 29] and start to train this model from scratch with both hard label and soft label supervision signals. The hard label denotes the ground-truth one-hot vector and the soft label represents the softened probability distribution from the original cumbersome model [8]. By distilling the probabilistic knowledge from the trained unpruned network, the small compact model can learn better representations and exhibit better generalization capability after the training phase. The final loss function is written as below, which is comprised of the hard label loss and the distillation loss:

$$\mathcal{L}_{retrain} = \mathcal{L}_{cls}(\mathbf{O}_S, \mathbf{O}_L) + \alpha_d \mathcal{L}_{distill}(\mathbf{P}_S, \mathbf{P}_T). \quad (6)$$

Here,  $\mathbf{O}_S$  is the probability output of the compact model,  $\mathbf{O}_L$  is the ground-truth label,  $\mathbf{P}_S$  and  $\mathbf{P}_T$  are the output

probabilities of the pruned and unpruned model, respectively. We denote  $\alpha_d$  as the coefficient to balance the hard label loss and distillation loss,  $\mathcal{L}_{cls}$  is the cross entropy loss and  $\mathcal{L}_{distill}$  is the KL divergence loss. Note that  $P_S$  and  $P_T$  are obtained via performing softmax on the output logits [8] and the temperature is set as 1.

## 4. Experiments

Following [3], we perform extensive experiments on ImageNet under various pruning levels. There are 1000 classes in ImageNet, around 128M images are used for training and 50K images for validation. We choose modern CNN architectures, *i.e.*, ResNet-18 [5] and ResNet-50 [5], for pruning as they are widely used in both academia and industry. Following [3], we also perform pruning on MobileNetV2 [23]. This is a more challenging task since MobileNetV2 is already a compact model. We also report results on pruning extremely efficient architectures, *i.e.*, MobileNetV3-small [12] and EfficientNet-B0 [26]. It is noteworthy that PEEL is not limited to these architectures and can be employed to address the pruning of more sophisticated networks.

**Evaluation metric:** Following [3], we adopt the top-1 accuracy on the center crop as the evaluation metric.

**Implementation details:** To ensure a fair comparison, we adopt the exact setting of [3]. Concretely, the total number of epochs for training the uniformly pruned backbone is 50. The final pruned architecture is trained for 100 epochs for ResNets, 250 epochs for MobileNetV2, and the first epoch is used to warm up the pruned model. Batch size is set as 512 and the initial learning rate is set as 0.2. The learning rate is gradually reduced to zero via a cosine learning rate schedule. Label smoothing is employed to relieve the over-fitting issue, and the corresponding smoothing factor is set as 0.1. Parameters  $\alpha_s$  and  $\alpha_d$  are set as 0.0001 and 0.1, respectively. As to data augmentation, random cropping and resizing, random horizontal flipping, color jitter and normalization are orderly applied to the input images during training. For knowledge distillation, ResNet-18 and ResNet-50 are trained with the original distillation strategy [8]. For MobileNetV2, we follow [3, 28] and make use of the in-place distillation to train the compact model.

**Baseline algorithms:** DMCP [3] is currently the most competitive algorithm in channel pruning. It models channel pruning as a Markov Decision Process, where the state is the remaining channels in a layer, and the transition from one state to another state represents the pruning process. By fusing the learnable transition probabilities with the layer output, the whole system is optimized end-to-end with gradient descent and can be easily combined with specific FLOPs constraint to acquire the desired network architecture. The training cost of DMCP is composed of two parts, *i.e.*, the learning of pruning probabilities and the training of

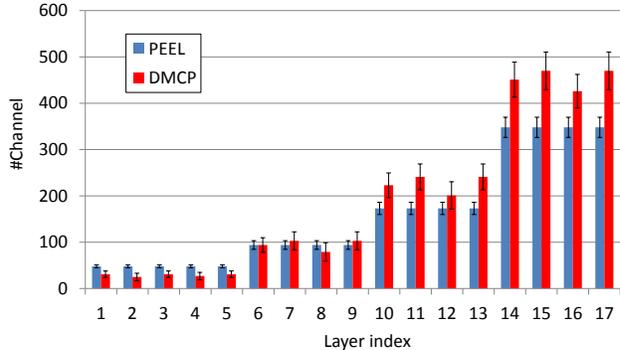


Figure 3: Comparison between structures searched by PEEL and DMCP on ResNet-18 when #FLOPs=1.04G. The vertical bar and vertical line denote the mean and variance of the channel numbers in each layer, respectively. Results are averaged over five runs.

the pruned model. The sampling procedure to acquire the compact model consumes negligible resources since they employ expected sampling to obtain the desired network in one round. In the subsequent paragraphs, we will systematically compare our PEEL with DMCP in terms of the overall performance under diverse FLOPs requirement, the expected training cost, the variability of the searched structures as well as performance stability. FPGM [7], NR [22], MetaPruning [17], AutoSlim [28] and AMC [6] are also included in the performance comparison since there are all contemporary channel pruning techniques.

### 4.1. Results

**Comparisons with state of the arts:** The detailed performance comparison of PEEL and contemporary channel pruning algorithms is shown in Table 1. We also record the floating-point operations (FLOPs) to approximate the forward time of different architectures. From Table. 1, we observe that PEEL achieves 25% and 45% FLOPs reduction for ResNet-18 and ResNet-50, respectively, with almost no loss of accuracy. As to MobileNetV2, the result is more encouraging as our algorithm can find a strong architecture with **2.5%** higher in top-1 accuracy than the original model while possessing the same FLOPs. In addition, PEEL consistently finds better architectures whose performance outperforms those searched by other competitive channel pruning techniques, *e.g.*, DMCP [3], NR [22] and MetaPruning [17]. For instance, when the backbone model is ResNet-50 and the target FLOPs is 2.2G, the architecture produced by PEEL is **0.5%** higher in top-1 accuracy than that found by DMCP. Moreover, the network structure obtained via PEEL consistently outperforms the uniformly pruned model in all experiments, suggesting the effectiveness of the proposed resource reallocation approach.

**Pruning MobileNetV3-small and EfficientNet-B0:** As

Table 1: Performance of PEEL and various channel pruning approaches with and without knowledge distillation (KD) on ImageNet validation set.

Backbone	Method	Top-1 Acc (%)		FLOPs	
		w/ KD	w/o KD		
ResNet-18	Uniform 1.0 ×	70.3		1.8G	
	Uniform 0.85 ×	69.0	68.5	1.33G	
	NR [22]	68.6	68.2	1.33G	
	DMCP [3]	70.0	69.7	1.33G	
	<b>PEEL</b>	<b>70.6</b>	<b>70.1</b>	<b>1.33G</b>	
	Uniform 0.75 ×	68.2	67.8	1.05G	
	NR [22]	68.1	67.9	1.04G	
	FPGM [7]	68.4	68.1	1.04G	
	DMCP [3]	69.4	69.0	1.04G	
	<b>PEEL</b>	<b>69.9</b>	<b>69.3</b>	<b>1.04G</b>	
	ResNet-50	Uniform 1.0 ×	76.4		4.1G
		Uniform 0.85 ×	75.8	75.4	3.0G
NR [22]		75.2	74.8	3.0G	
MetaPruning [17]		76.2	76.0	3.0G	
AutoSlim [28]		76.0	75.8	3.0G	
DMCP [3]		76.5	76.4	3.0G	
<b>PEEL</b>		<b>76.9</b>	<b>76.6</b>	<b>3.0G</b>	
Uniform 0.75 ×		74.6	74.1	2.3G	
NR [22]		74.3	74.0	2.4G	
MetaPruning [17]		75.4	75.2	2.3G	
FPGM [7]		75.6	75.5	2.4G	
DMCP [3]		76.3	76.2	2.2G	
<b>PEEL</b>		<b>76.8</b>	<b>76.5</b>	<b>2.2G</b>	
Uniform 0.5 ×		72.9	72.7	1.1G	
NR [22]		73.1	72.6	1.1G	
MetaPruning [17]		73.4	73.2	1.1G	
AutoSlim [28]		74.0	73.6	1.1G	
DMCP [3]		74.6	74.4	1.1G	
<b>PEEL</b>		<b>75.1</b>	<b>74.7</b>	<b>1.1G</b>	
MobileNetV2		Uniform 1.0 ×	72.3		300M
		NR [22]	73.0	72.2	300M
		AutoSlim [28]	74.2	73.1	300M
		DMCP [3]	74.6	73.9	300M
		<b>PEEL</b>	<b>74.8</b>	<b>74.2</b>	<b>300M</b>
	Uniform 0.75 ×	70.1	69.2	210M	
	NR [22]	70.2	69.4	211M	
	MetaPruning [17]	71.2	70.1	217M	
	AMC [6]	70.8	70.0	211M	
	AutoSlim [28]	73.0	72.2	211M	
	DMCP [3]	73.5	72.4	211M	
	<b>PEEL</b>	<b>73.9</b>	<b>73.0</b>	<b>211M</b>	
	Uniform 0.5 ×	64.8	64.3	97M	
	NR [22]	64.2	63.8	87M	
	MetaPruning [17]	63.8	63.4	87M	
	DMCP [3]	66.1	65.6	87M	
	<b>PEEL</b>	<b>66.6</b>	<b>66.0</b>	<b>87M</b>	
	Uniform 0.35 ×	60.1	59.7	59M	
	NR [22]	59.7	59.2	59M	
	DMCP [3]	62.7	62.4	59M	
	<b>PEEL</b>	<b>62.9</b>	<b>62.6</b>	<b>59M</b>	

Table 2: Pruning MobileNetV3-small and EfficientNet-B0.

Backbone	Method	Top-1 Acc (%)	FLOPs
MobileNetV3	Uniform 1.0 ×	67.4	56M
	<b>PEEL</b>	<b>67.2</b>	<b>49M</b>
EfficientNet-B0	Uniform 1.0 ×	77.1	390M
	<b>PEEL</b>	<b>77.0</b>	<b>346M</b>

Table 3: Performance variance of searched architectures of PEEL and DMCP on ResNet-50. Results are averaged over five runs.

Method	#FLOPs		
	2.2G	1.8G	1.1G
<b>PEEL</b>	<b>76.7% ± 0.2%</b>	<b>75.3% ± 0.3%</b>	<b>75.0% ± 0.2%</b>
DMCP	75.7% ± 0.7%	74.2% ± 0.8%	73.8% ± 0.8%

Table 4: Comparison between the training time and memory usage of PEEL and DMCP on ResNet-50 (#FLOPs=1.1G). Here, the memory usage is measured on one GPU and the batch size is set as 64 for each GPU.

Algorithm	Train Time (H)	Memory usage (G)
<b>PEEL</b>	<b>37</b>	<b>4.3</b>
DMCP [3]	98	9.5

shown in Table 2, PEEL saves **12.5%** FLOPs for MobileNetV3-small and **11.3%** FLOPs for EfficientNet-B0 without incurring severe performance drops. Note that MobileNetV3-small and EfficientNet-B0 are two extremely lightweight classification models and removing redundancy in these models is nontrivial. The experimental results explicitly showcase the effectiveness and generality of PEEL on network pruning.

We also visualize the structures uncovered by PEEL and DMCP. As depicted in Fig. 3, both architectures have more channels as the layer goes deeper. This is expected since more channels will be leveraged to compensate for the loss of spatial information incurred by the downsampling operations. Nevertheless, PEEL tends to place more channels in the early stages while DMCP chooses to put more channels in the later stages. Putting more resources at shallow layers brings more gains as there is more redundancy in the deep layers. Besides, the variance of the number of channels in each layer of PEEL is much smaller than that of DMCP. For example, the averaged channel variance of DMCP on the last four layers of ResNet-18 is **38.8** while the channel variance of PEEL is **21.8**. The gap between the structural variance of DMCP and our algorithm is more evident in ResNet-50 and MobileNetV2. In ResNet-50, the averaged channel variance of PEEL is almost half of the variance of DMCP for all layers. We provide more details in the supplementary material. The searching instability of DMCP also leads to high performance variance in the searched architectures. The performance of PEEL is much stabler relatively. We summarize the performance of architectures searched by PEEL and DMCP in Table 3.

The proposed PEEL not only takes shorter training durations to find and train the desired model, but also consumes less GPU memory during the searching phase. Here, we examine the training cost of PEEL and DMCP, *i.e.*, the total training hours and the GPU memory usage during the searching phase. The original unpruned model is ResNet-

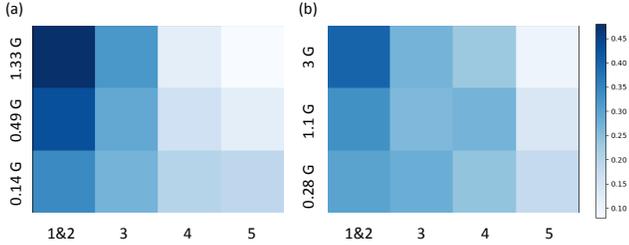


Figure 4: Where the resources go. Percentage of FLOPs assigned by PEEL to different groups of (a) ResNet-18 and (b) ResNet-50 under different target FLOPs constraints. X-axis denotes group indices and y-axis represents the target FLOPs. A darker color indicates a larger quantity of assigned resources. Since ResNet-18 and ResNet-50 have only one layer in group 1 and resources assigned to one layer are very limited, we put this layer into group 2 and obtain the above figure.

50 and the target FLOPs is 1.1 G. We keep all other hyperparameters the same, including batch size (512), GPU types (NVIDIA TITAN X) and number of used GPUs (8). From Table 4, the total number of training hours of PEEL is roughly a third of DMCP. The result is not surprising as DMCP requires training on the original cumbersome model whilst PEEL merely entails training of a much slimmer model. As to the GPU memory usage, since DMCP needs to collect gradients of several sampled architectures while PEEL only computes the gradient of one compact architecture, the memory consumption of DMCP is almost twice as that of PEEL. We also provide a comparison to another representative USNet [29] in the supplementary material.

## 4.2. Ablation study

**Where the resources go:** To have a deeper understanding on the effect of the resource reallocation module, we visualize the percentage of resources distributed by PEEL to different groups of layers on ResNet-18 and ResNet-50. As shown in Fig. 4, layers in the first three groups are given much more resources than the last two groups when the network is slightly trimmed. It is natural as there are more redundant channels in the deep layers in the mild pruning level [17]. As the pruning becomes more aggressive, the percentages of FLOPs assigned to different groups become more even since all layers do not have sufficient channels and call for more resources from the resource pool. And we can observe the same patterns in the resource reallocation of ResNet-18 and ResNet-50 when the FLOPs budget becomes tighter.

**Effect of  $\lambda$ :** The resources available in the resource pool is  $(1 - \lambda)M$ , thus the pool size is controlled by the hyperparameter  $\lambda$ . We select the value of  $\lambda$  from  $\{0.5, 0.6, 0.7, 0.8, 0.9, 1\}$  and compare the model performance under these settings. As illustrated in Fig. 5 (a),

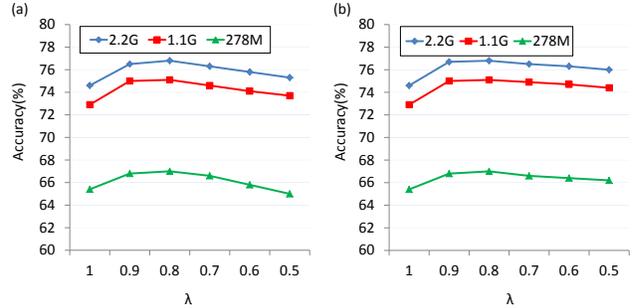


Figure 5: One-round v.s. multi-round reallocation. We show the  $\lambda$ -accuracy spectrum of (a) one-round and (b) three-round resource reallocation of PEEL on ResNet-50. Different color denotes different target FLOPs.

the performance of the searched model is relatively stable when  $\lambda$  ranges from 0.7 to 0.9. When  $\lambda$  decreases from 0.7, the performance of the final pruned architecture gradually drops. The trend is expected since PEEL only performs the evaluation of layer importance once on the over-pruned model. The importance evaluation would become inaccurate given a small  $\lambda$ , as more resources are assigned to the pool while the over-pruned backbone is too slim. Another phenomenon we observe is that as the pruning becomes more aggressive, the performance of the pruned model is more susceptible to the value of  $\lambda$  (see the green line in Fig. 5 (a)). We conjecture that the performance drop is caused by the challenging pruning condition. Given very limited FLOPs resources, one has to allocate resources in a very careful manner so that the pruned model can exhibit satisfactory performance.

### One-round reallocation v.s. multi-round reallocation:

Recall that we adopt the one-round resource reallocation in our algorithm, *i.e.*, allocating the resources by just performing a single pass of layer importance estimation. An alternative strategy is to evenly divide the resources into several parts and reallocate these resources successively in multiple rounds. Multi-round reallocation is more expensive as one needs to repeat Step-2 and Step-4 iteratively (see Sec. 3), and each iteration involves the fine-tuning of the backbone. As depicted in Fig. 5, the performance of the multi-round resource reallocation is more stable than the one-round version when the  $\lambda$  is set 0.6 and 0.5. For instance, when the target FLOPs is 2.2G and  $\lambda$  is set as 0.5, the performance of the one-round reallocation decreases from 76.8% to 75.3% while the multi-round reallocation can still achieve 76.0% in top-1 accuracy. The more stable performance may come from the fact that multi-round resource reallocation reevaluates the importance of each layer in each round, thus allowing a more appropriate resource reallocation.

**Robustness to the layer importance indicator:** The original PEEL adopts BN statistics to reflect the importance of each layer. Here, we explore different importance in-

Table 5: Performance of PEEL with different importance criteria on ResNet-50.

Criterion	#FLOPs		
	2.2 G	1.8 G	1.1 G
BN statistics (Ours)	76.8%	75.5%	75.1%
Filter norm [14]	76.7%	75.2%	75.2%
Reconstruction errors [20]	76.6%	75.4%	75.0%

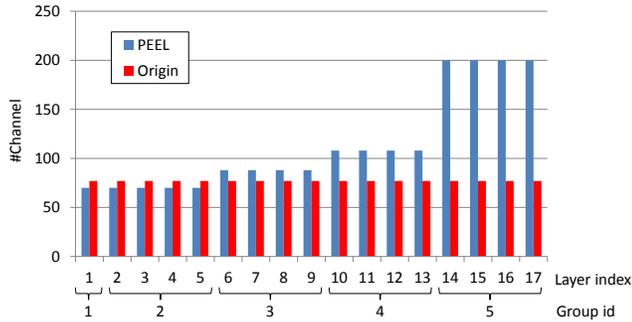


Figure 6: Visualization of the architecture found by PEEL on a different backbone model that has the same number of channels for all layers. The base model is ResNet-18, and target FLOPs is 1.04G.

indicators and check whether the proposed PEEL is sensitive to the chosen indicator. We choose the widely-used filter norm [14] and reconstruction errors [20] as the criteria. Layers with large filter norm or large reconstruction errors are considered as important. However, directly taking filter norm as the evaluating metric is biased since the filter norm of different layers have diverse scales in magnitude [1]. To correct such bias, we follow [1] and learn the layer-wise affine transformations for the filter norm. Under this circumstance, filter norm can better reflect the value of each layer. From Table 5, when pruning ResNet-50 under different FLOPs requirements, architectures searched by PEEL using BN statistics, filter norm and reconstruction errors achieve similar performance. The results suggest the robustness of PEEL, irrespective of the methods used in evaluating layer importance.

**Effect of knowledge distillation:** By comparing the two columns of the Top-1 Acc in Table 1, we find that PEEL consistently yields better performance than baseline pruning approaches with or without knowledge distillation (KD). For instance, when performing pruning on ResNet-18 and target FLOPs is 1.33G, PEEL with KD is 0.6% higher than DMCP with KD. One interesting phenomenon we observe is that KD brings more gains when the backbone model has fewer parameters. For example, KD can bring approximately 0.3%, 0.5% and 1.0% to PEEL on ResNet-50, ResNet-18 and MobileNetV2, respectively when the pruning is not aggressive. We conjecture that the increased gains are attributed to the deficiency of small models in

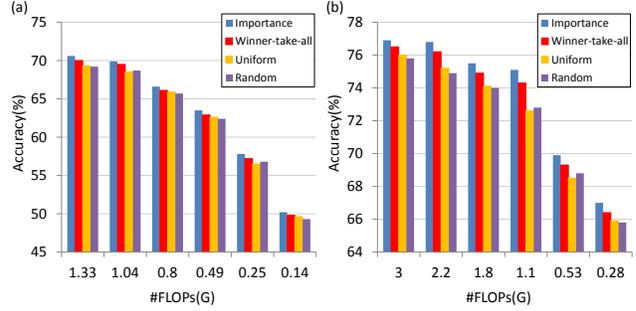


Figure 7: Policy on resource reallocation. Performance of PEEL with importance-guided, winner-take-all, uniform and random policies on (a) ResNet-18 and (b) ResNet-50, respectively.

grasping knowledge in labels by themselves and hence they more eagerly call for the guidance of the original model.

**Efficacy of PEEL on a different backbone model:** Recall that we treat the uniformly trimmed model as the backbone and conduct resource reallocation on it. It is natural to wonder whether the resource reallocation would still work if a totally distinct architecture is used. Here, we use an architecture that has the same number of channels across all layers as the new starting point. All configurations are the same as the previous experiments. From Fig. 6, we can observe that in this new backbone, groups 1, 2 and 3 are assigned with few parameters while the majority of the resources are distributed to groups 4 and 5. Overall, PEEL puts more parameters in the deeper layers, and the resulting architecture is 8.7% higher in terms of top-1 accuracy in comparison to the original model (62.4% v.s. 53.7%). The result demonstrates the effectiveness of resource reallocation as well as the insensitivity of our algorithm to the backbone model.

**Policy on resource reallocation:** Here, we compare our ‘importance-guided’ resource reallocation strategy with other parameter reallocation policies, *i.e.*, winner-take-all, uniform and random reassignment. The winner-take-all policy determines the most important group according to the computed group significance and puts all resources in that group. The uniform policy adds the same number of channels to all layers while the random policy stochastically selects several layers and increases their channels. Here, we name the original reallocation strategy as importance-guided policy as it distributes resources based upon the estimated layer importance. From Table 7, our importance-guided reallocation policy evidently outperforms the other three policies on ResNet-18 and ResNet-50. For instance, when the target FLOPs is 1.8G on ResNet-50, the top-1 accuracy of the importance-guided policy is 75.5% while the top-1 accuracy of the winner-take-all, uniform and random policy is 74.9%, 74.1% and 74.0%, respectively. The superior performance demonstrates the effectiveness of the importance-guided policy.

## 5. Conclusion

We have presented an easy-to-implement yet effective channel pruning algorithm, *i.e.*, PEEL, to acquire a desired compact model via reallocating resources from less informative layers to more crucial layers. The intuition of PEEL is that layers do not contribute equally to the model performance and those having an indispensable influence on the performance should be assigned more resources to magnify their positive effect. We conduct extensive experiments to verify the efficacy of PEEL on ImageNet dataset with modern CNN architectures, *i.e.*, ResNet-18, ResNet-50, MobileNetV2, MobileNetV3-small and EfficientNet-B0. Experimental results suggest the effectiveness of PEEL in uncovering compact yet accurate architectures consistently under various pruning levels, compared with state-of-the-art channel pruning methods. Besides, PEEL yields stable searching results with small performance variance, and the searching cost is evidently smaller than contemporary channel pruning algorithms (*e.g.*, DMCP [3]).

## References

- [1] Ting-Wu Chin, Ruizhou Ding, Cha Zhang, and Diana Marculescu. Towards Efficient Model Compression via Learned Global Ranking. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1518–1528, 2020. 8
- [2] Wei Gao, Yi Wei, Quanquan Li, Hongwei Qin, Wanli Ouyang, and Junjie Yan. Pruning with Hints: An Efficient Framework for Model Acceleration. 2018. 2
- [3] Shaopeng Guo, Yujie Wang, Quanquan Li, and Junjie Yan. DMCP: Differentiable Markov Channel Pruning for Neural Networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1539–1547, 2020. 1, 2, 3, 4, 5, 6, 9, 10, 11
- [4] Song Han, Jeff Pool, John Tran, and William Dally. Learning Both Weights and Connections for Efficient Neural Network. In *Advances in Neural Information Processing Systems*, pages 1135–1143, 2015. 2, 4
- [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 1, 5, 11
- [6] Yihui He, Ji Lin, Zhijian Liu, Hanrui Wang, Li-Jia Li, and Song Han. AMC: Automl for Model Compression and Acceleration on Mobile Devices. In *European Conference on Computer Vision (ECCV)*, pages 784–800, 2018. 2, 5, 6
- [7] Yang He, Ping Liu, Ziwei Wang, Zhilan Hu, and Yi Yang. Filter Pruning via Geometric Median for Deep Convolutional Neural Networks Acceleration. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4340–4349, 2019. 2, 5, 6
- [8] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the Knowledge in a Neural Network. *Statistics*, 1050:9, 2015. 3, 4, 5
- [9] Yuenan Hou, Zheng Ma, Chunxiao Liu, Tak-Wai Hui, and Chen Change Loy. Inter-Region Affinity Distillation for Road Marking Segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 12486–12495, 2020. 3
- [10] Yuenan Hou, Zheng Ma, Chunxiao Liu, and Chen Change Loy. Learning Lightweight Lane Detection CNNs by Self Attention Distillation. In *IEEE International Conference on Computer Vision*, pages 1013–1021, 2019. 3
- [11] Yuenan Hou, Zheng Ma, Chunxiao Liu, and Chen Change Loy. Learning to Steer by Mimicking Features from Heterogeneous Auxiliary Networks. In *Association for the Advancement of Artificial Intelligence*, volume 33, pages 8433–8440, 2019. 3
- [12] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, Quoc V. Le, and Hartwig Adam. Searching for mobilenetv3. In *IEEE International Conference on Computer Vision*, October 2019. 5
- [13] Sergey Ioffe and Christian Szegedy. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. *International Conference on Machine Learning*, 2015. 4
- [14] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning Filters for Efficient Convnets. In *International Conference on Learning Representations*, 2017. 1, 8
- [15] Yifan Liu, Ke Chen, Chris Liu, Zengchang Qin, Zhenbo Luo, and Jingdong Wang. Structured Knowledge Distillation for Semantic Segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2604–2613, 2019. 3
- [16] Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang. Learning Efficient Convolutional Networks through Network Slimming. In *IEEE International Conference on Computer Vision*, pages 2736–2744, 2017. 1, 2, 3, 4
- [17] Zechun Liu, Haoyuan Mu, Xiangyu Zhang, Zichao Guo, Xin Yang, Kwang-Ting Cheng, and Jian Sun. Metapruning: Meta Learning for Automatic Neural Network Channel Pruning. In *IEEE International Conference on Computer Vision*, pages 3296–3305, 2019. 3, 5, 6, 7
- [18] Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. Rethinking the Value of Network Pruning. In *International Conference on Learning Representations*, 2018. 3
- [19] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully Convolutional Networks for Semantic Segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015. 1
- [20] Jian-Hao Luo, Jianxin Wu, and Weiyao Lin. Thinet: A Filter Level Pruning Method for Deep Neural Network Compression. In *IEEE International Conference on Computer Vision*, pages 5058–5066, 2017. 1, 2, 4, 8
- [21] Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning Convolutional Neural Networks for Resource Efficient Inference. In *International Conference on Learning Representations*, 2017. 2, 4
- [22] Siyuan Qiao, Zhe Lin, Jianming Zhang, and Alan L Yuille. Neural Rejuvenation: Improving Deep Network Training by

- Enhancing Computational Resource Utilization. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 61–71, 2019. 2, 5, 6
- [23] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. MobileNetV2: Inverted Residuals and Linear Bottlenecks. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4510–4520, 2018. 5, 10, 11
- [24] Dawei Sun, Anbang Yao, Aojun Zhou, and Hao Zhao. Deeply-Supervised Knowledge Synergy. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 6997–7006, 2019. 3
- [25] Christian Szegedy, Alexander Toshev, and Dumitru Erhan. Deep Neural Networks for Object Detection. In *Advances in Neural Information Processing Systems*, pages 2553–2561, 2013. 1
- [26] Mingxing Tan and Quoc V. Le. EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. *International Conference on Machine Learning*, 2019. 5
- [27] Mao Ye, Chengyue Gong, Lizhen Nie, Denny Zhou, Adam Klivans, and Qiang Liu. Good Subnetworks Provably Exist: Pruning via Greedy Forward Selection. *International Conference on Machine Learning*, 2020. 2
- [28] Jiahui Yu and Thomas Huang. Autoslim: Towards One-Shot Architecture Search for Channel Numbers. *arXiv preprint arXiv:1903.11728*, 2019. 2, 5, 6
- [29] Jiahui Yu and Thomas S Huang. Universally Slimmable Networks and Improved Training Techniques. In *IEEE International Conference on Computer Vision*, pages 1803–1811, 2019. 1, 2, 3, 4, 7, 10, 11
- [30] Sergey Zagoruyko and Nikos Komodakis. Paying More Attention to Attention: Improving the Performance of Convolutional Neural Networks via Attention Transfer. In *International Conference on Learning Representations*, 2017. 3
- [31] Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be Your Own Teacher: Improve the Performance of Convolutional Neural Networks via Self Distillation. In *IEEE International Conference on Computer Vision*, pages 3713–3722, 2019. 3
- [32] Zhuangwei Zhuang, Mingkui Tan, Bohan Zhuang, Jing Liu, Yong Guo, Qingyao Wu, Junzhou Huang, and Jinhui Zhu. Discrimination-Aware Channel Pruning for Deep Neural Networks. In *Advances in Neural Information Processing Systems*, pages 875–886, 2018. 2

## Appendices

### A. A detailed comparison with DMCP

We compare PEEL against the recent and representative DMCP [3]. Figure 8 shows the accuracy of architectures pruned by PEEL and DMCP under diverse FLOPs constraints on ResNet-18, ResNet-50 and MobileNetV2. Here, we take the uniform pruning as reference. It is observed that PEEL can always find better architectures than DMCP

Table 6: Performance variance of searched architectures of PEEL and DMCP on ResNet-50. Results are averaged over five runs.

Method	#FLOPs		
	2.2G	1.8G	1.1G
<b>PEEL</b>	<b>76.7% ± 0.2%</b>	<b>75.3% ± 0.3%</b>	<b>75.0% ± 0.2%</b>
DMCP [3]	75.7% ± 0.7%	74.2% ± 0.8%	73.8% ± 0.8%
USNet [29]	75.4% ± 0.6%	74.0% ± 0.7%	73.4% ± 0.6%

Table 7: Comparison between the training time and memory usage of PEEL and DMCP on ResNet-50 (#FLOPs=1.1G). Here, the memory usage is measured on one GPU and the batch size is set as 64 for each GPU.

Algorithm	Train Time (H)	Memory usage (G)
<b>PEEL</b>	<b>37</b>	<b>4.3</b>
DMCP [3]	98	9.5
USNet [29]	66	8.7

in terms of top-1 accuracy. On all FLOPs levels, PEEL outperforms DMCP by least 0.4% in terms of top-1 accuracy. Thanks to the resource reallocation, the top-1 accuracy gap between PEEL and uniform pruning lies between 1.0% and 3.8%.

Figure 9 reveals the detailed FLOPs assignment on MobileNetV2 [23]. We observe a similar pattern in the FLOPs distribution process of MobileNetV2 to that of ResNet-18 and ResNet-50. Concretely, when the pruning is mild, the resource pool will tend to allocate more resources in the first three groups as more redundancy is expected in deep layers. When the target FLOPs budget decreases, the reallocated resources will be more even for five groups as all groups are in need of more parameters to function effectively.

### B. $\lambda$ -accuracy spectrum of one-round and multi-round resource reallocation

We have shown the results of ResNet-50 with one-round and multi-round reallocation in Figure 6 of the main paper. And here, we provide detailed results on ResNet-18 and MobileNetV2. It is noteworthy that we observe a similar trend on ResNet-18 and MobileNetV2. Specifically, as shown in Fig. 10 (1), PEEL with one-round reallocation can produce relatively stable results when  $\lambda$  is not smaller than 0.7. When  $\lambda$  continues to decrease, the performance of PEEL suffers as the one-round evaluation of layer importance becomes less accurate. PEEL with multi-round resource reallocation (Fig. 10 (2)) can avoid severe performance drop when  $\lambda$  is set as 0.6 or 0.5 since it performs the importance assessment during each round, thus yielding more reliable statistics of layer importance.

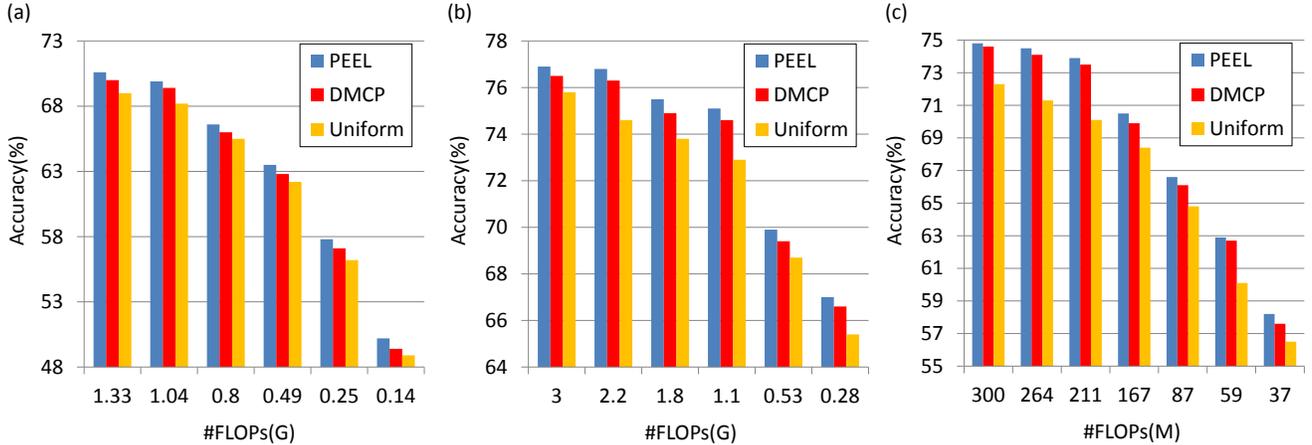


Figure 8: FLOPs-accuracy spectrum of PEEL, DMCP and uniform pruning on (a) ResNet-18, (b) ResNet-50 and (c) MobileNetV2.

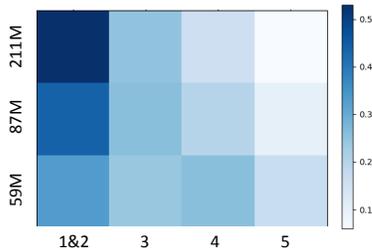


Figure 9: Percentage of FLOPs assigned by PEEL to different groups of MobileNetV2 under different target FLOPs constraints. X-axis denotes group indices and y-axis represents the target FLOPs. A darker color indicates a larger quantity of assigned resources.

### C. Visualization of architectures searched by PEEL and DMCP

We visualize the architectures obtained by PEEL and DMCP [3] on ResNet-50 [5] and MobileNetV2 [23]. As depicted in Fig. 11, on both ResNet-50 and MobileNetV2, PEEL puts more resources on shallow layers whilst DMCP has more reallocated parameters on the deep layers. As to the architectural stability, the variance of channel numbers of PEEL is much smaller than that of DMCP. For instance, on ResNet-50, the channel variance of the last nine layers of PEEL is less than half of DMCP (28.8 v.s. 64.3). These results strongly support the searching stability of PEEL.

### D. Comparison between performance variance and searching cost of PEEL, DMCP and USNet

We summarize the performance variance and searching cost of different pruning algorithms in Table. 6 and Table. 7, respectively. It is evident that PEEL has less

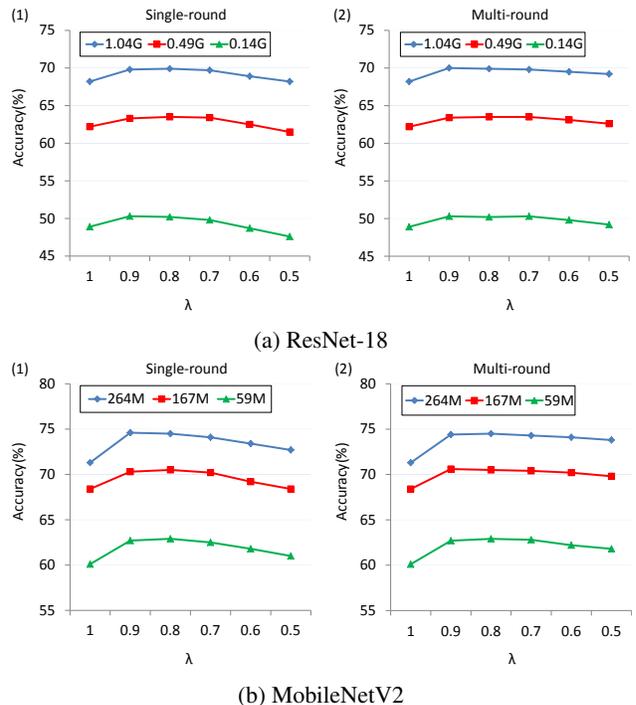
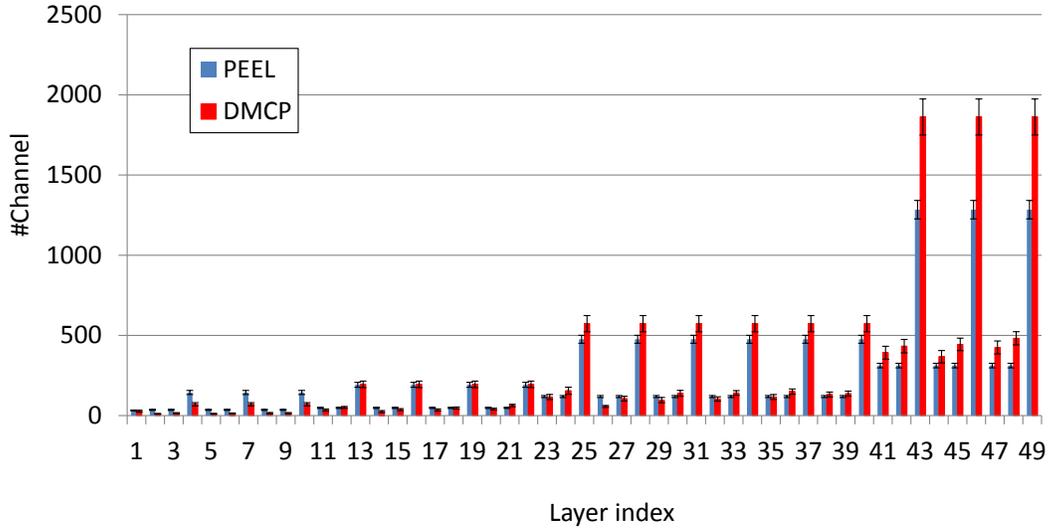
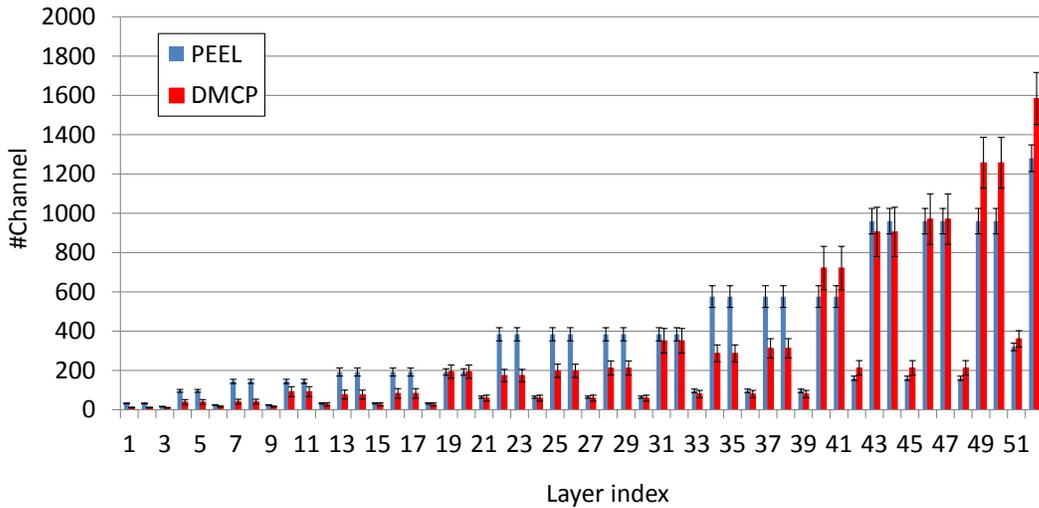


Figure 10: One-round v.s. multi-round reallocation. We show the  $\lambda$ -accuracy spectrum of one-round (left) and multi-round (right) resource reallocation of PEEL on (a) ResNet-18 and (b) MobileNetV2, respectively. Different color denotes different target FLOPs.

variance in performance than DMCP [3] and USNet [29]. Besides, PEEL takes shorter training hours and consumes much less GPU memory to produce a desired slim model than the other two pruning techniques. It is not surprising since PEEL trains a much compact model and reallocates resources on this model while both DMCP and USNet en-



(a) ResNet-50



(b) MobileNetV2

Figure 11: Comparison between structures searched by PEEL and DMCP on (a) ResNet-50 when #FLOPs=1.1G and (b) MobileNetV2 when #FLOPs=300M, respectively. The vertical bar and vertical line denote the mean and variance of the channel numbers in each layer, respectively. Results are averaged over five runs.

tails the training of the original cumbersome network. Besides, as opposed to DMCP and USNet, PEEL is free from calculating the gradients of multiple sampled substructures, thus occupying fewer memory resources during the searching phase. These results explicitly showcases the superior efficiency of PEEL.