

Towards Efficient Model Compression via Learned Global Ranking

Ting-Wu Chin¹, Ruizhou Ding¹, Cha Zhang², Diana Marculescu^{1,3}
Carnegie Mellon University¹, Microsoft Cloud and AI², The University of Texas at Austin³
{tingwuc, rding}@andrew.cmu.edu, chazhang@microsoft.com, dianam@utexas.edu

Abstract

*Pruning convolutional filters has demonstrated its effectiveness in compressing ConvNets. Prior art in filter pruning requires users to specify a target model complexity (e.g., model size or FLOP count) for the resulting architecture. However, determining a target model complexity can be difficult for optimizing various embodied AI applications such as autonomous robots, drones, and user-facing applications. First, both the accuracy and the speed of ConvNets can affect the performance of the application. Second, the performance of the application can be hard to assess without evaluating ConvNets during inference. As a consequence, finding a sweet-spot between the accuracy and speed via filter pruning, which needs to be done in a trial-and-error fashion, can be time-consuming. This work takes a first step toward making this process more efficient by altering the goal of model compression to producing a set of ConvNets with various accuracy and latency trade-offs instead of producing one ConvNet targeting some pre-defined latency constraint. To this end, we propose to learn a **global ranking** of the filters across different layers of the ConvNet, which is used to obtain a set of ConvNet architectures that have different accuracy/latency trade-offs by pruning the bottom-ranked filters. Our proposed algorithm, LeGR, is shown to be 2× to 3× faster than prior work while having comparable or better performance when targeting seven pruned ResNet-56 with different accuracy/FLOPs profiles on the CIFAR-100 dataset. Additionally, we have evaluated LeGR on ImageNet and Bird-200 with ResNet-50 and MobileNetV2 to demonstrate its effectiveness. Code available at <https://github.com/cmu-enyac/LeGR>.*

1. Introduction

Building on top of the success of visual perception [49, 17, 18], natural language processing [10, 11], and speech recognition [6, 45] with deep learning, researchers have started to explore the possibility of embodied AI applications. In embodied AI, the goal is to enable agents to take actions based on perceptions in some environments [51].

We envision that next generation embodied AI systems will run on mobile devices such as autonomous robots and drones, where compute resources are limited and thus, will require model compression techniques for bringing such intelligent agents into our lives.

In particular, pruning the convolutional filters in ConvNets, also known as filter pruning, has shown to be an effective technique [63, 36, 60, 32] for trading accuracy for inference speed improvements. The core idea of filter pruning is to find the least important filters to prune by minimizing the accuracy degradation and maximizing the speed improvement. State-of-the-art filter pruning methods [16, 20, 36, 70, 46, 8] require a target model complexity of the whole ConvNet (e.g., total filter count, FLOP count¹, model size, inference latency, etc.) to obtain a pruned network. However, deciding a target model complexity for optimizing embodied AI applications can be hard. For example, considering delivery with autonomous drones, both inference speed and precision of object detectors can affect the drone velocity [3], which in turn affects the inference speed and precision². For an user-facing autonomous robot that has to perform complicated tasks such as MovieQA [56], VQA [2], and room-to-room navigation [1], both speed and accuracy of the visual perception module can affect the user experience. These aforementioned applications require many iterations of trial-and-error to find the optimal trade-off point between speed and accuracy of the ConvNets.

More concretely, in these scenarios, practitioners would have to determine the sweet-spot for model complexity and accuracy in a trial-and-error fashion. Using an existing filter pruning algorithm many times to explore the impact of the different accuracy-vs.-speed trade-offs can be time-consuming. Figure 1 demonstrates the usage of filter pruning for optimizing ConvNets in aforementioned scenarios. With prior approaches, one has to go through the process of finding constraint-satisfying pruned-ConvNets via a pruned-

¹The number of floating-point operations to be computed for a ConvNet to carry out an inference.

²Higher velocity requires faster computation and might cause accuracy degradation due to the blurring effect of the input video stream.

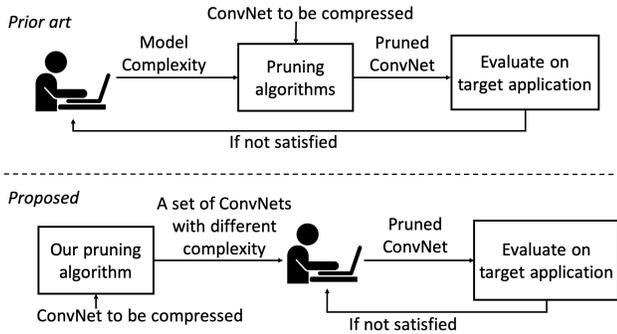


Figure 1: Using filter pruning to optimize ConvNets for embodied AI applications. Instead of producing one ConvNet for each pruning procedure as in prior art, our proposed method produces a set of ConvNets for practitioners to efficiently explore the trade-offs.

ing algorithm for every model complexity considered until practitioners are satisfied with the accuracy-vs.-speedup trade-off. Our work takes a first step toward alleviating the inefficiency in the aforementioned paradigm. We propose to alter the objective of pruning from outputting a single ConvNet with pre-defined model complexity to producing a set of ConvNets that have different accuracy/speed trade-offs, while achieving comparable accuracy with state-of-the-art methods (as shown in Figure 4). In this fashion, the model compression overhead can be greatly reduced, which results in a more practical usage of filter pruning.

To this end, we propose *learned global ranking (or LeGR)*, an algorithm that learns to rank convolutional filters across layers such that the ConvNet architectures of different speed/accuracy trade-offs can be obtained easily by dropping the bottom-ranked filters. The obtained architectures are then fine-tuned to generate the final models. In such a formulation, one can obtain a set of architectures by learning the ranking *once*. We demonstrate the effectiveness of the proposed method with extensive empirical analyses using ResNet and MobileNetV2 on CIFAR-10/100, Bird-200, and ImageNet datasets. The main contributions of this work are as follows:

- We propose learned global ranking (*LeGR*), which produces a set of pruned ConvNets with different accuracy/speed trade-offs. *LeGR* is shown to be faster than prior art in ConvNet pruning, while achieving comparable accuracy with state-of-the-art methods on three datasets and two types of ConvNets.
- Our formulation towards pruning is the first work that considers learning to rank filters across different layers globally, which addresses the limitation of prior art in magnitude-based filter pruning.

2. Related Work

Various methods have been developed to compress and/or accelerate ConvNets including weight quantization [47, 71, 29, 30, 65, 24, 12, 7], efficient convolution operators [25, 22, 61, 26, 67], neural architecture search [69, 9, 4, 15, 54, 53, 52], adjusting image resolution [55, 5], and filter pruning, considered in this paper. Prior art on filter pruning can be grouped into two classes, depending on whether the architecture of the pruned-ConvNet is assumed to be given.

Pre-defined architecture In this category, various work proposes different metrics to evaluate the importance of filters locally within each layer. For example, some prior work [32, 19] proposes to use ℓ_2 -norm of filter weights as the importance measure. On the other hand, other work has also investigated using the output discrepancy between the pruned and unpruned network as an importance measure [23, 40]. However, the key drawback for methods that rank filters locally within a layer is that it is often hard to decide the overall target pruned architectures [20]. To cope with this difficulty, uniformly pruning the same portion of filters across all the layers is often adopted [19].

Learned architecture In this category, pruning algorithms learn the resulting structure automatically given a controllable parameter to determine the complexity of the pruned-ConvNet. To encourage weights with small magnitudes, Wen *et al.* [60] propose to add group-Lasso regularization to the filter norm to encourage filter weights to be zeros. Later, Liu *et al.* [36] propose to add Lasso regularization on the *batch normalization layer* to achieve pruning during training. Gordon *et al.* [16] propose to add compute-weighted Lasso regularization on the filter norm. Huang *et al.* [27] propose to add Lasso regularization on the output neurons instead of weights. While the regularization pushes unimportant filters to have smaller weights, the final thresholding applied globally assumes different layers to be equally important. Later, Louizos *et al.* [39] have proposed L_0 regularization with stochastic relaxation. From a Bayesian perspective, Louizos *et al.* [38] formulate pruning in a probabilistic fashion with a sparsity-induced prior. Similarly, Zhou *et al.* [70] propose to model inter-layer dependency. From a different perspective, He *et al.* propose an automated model compression framework (AMC) [20], which uses reinforcement learning to search for a ConvNet that satisfies user-specified complexity constraints.

While these prior approaches provide competitive pruned-ConvNets under a given target model complexity, it is often hard for one to specify the complexity parameter when compressing a ConvNet in embodied AI applications. To cope with this, our work proposes to generate a set of

pruned-ConvNets across different complexity values rather than a single pruned-ConvNet under a target model complexity.

We note that some prior work gradually prunes the ConvNet by alternating between pruning out a filter and fine-tuning, and thus, can also obtain a set of pruned-ConvNets with different complexities. For example, Molchanov *et al.* [43] propose to use the normalized Taylor approximation of the loss as a measure to prune filters. Specifically, they greedily prune one filter at a time and fine-tune the network for a few gradient steps before the pruning proceeds. Following this paradigm, Theis *et al.* [57] propose to switch from first-order Taylor to Fisher information. However, our experiment results show that the pruned-ConvNet obtained by these methods have inferior accuracy compared to the methods that generate a single pruned ConvNet.

To obtain a set of ConvNets across different complexities with competitive performance, we propose to learn a global ranking of filters across different layers in a data-driven fashion such that architectures with different complexities can be obtained by pruning out the bottom-ranked filters.

3. Learned Global Ranking

The core idea of the proposed method is to learn a ranking for filters across different layers such that a ConvNet of a given complexity can be obtained easily by pruning out the bottom rank filters. In this section, we discuss our assumptions and formulation toward achieving this goal.

As mentioned earlier in Section 1, often both accuracy and latency of a ConvNet affect the performance of the overall application. The goal for model compression in these settings is to explore the accuracy-vs.-speed trade-off for finding a sweet-spot for a particular application using model compression. Thus, in this work, we use FLOP count for the model complexity to sample ConvNets. As we will show in Section 5.3, we find FLOP count to be predictive for latency.

3.1. Global Ranking

To obtain pruned-ConvNets with different FLOP counts, we propose to learn the filter ranking globally across layers. In such a formulation, the global ranking for a given ConvNet just needs to be learned once and can be used to obtain ConvNets with different FLOP counts. However, there are two challenges for such a formulation. First, the global ranking formulation enforces an assumption that the top-performing smaller ConvNets are a proper subset of the top-performing larger ConvNets. The assumption might be strong because there are many ways to set the filter counts across different layers to achieve a given FLOP count, which implies that there are opportunities where the top-performing smaller network can have more filter counts in some layers but fewer filter counts in some other layers

compared to a top-performing larger ConvNet. Nonetheless, this assumption enables the idea of global filter ranking, which can generate pruned ConvNets with different FLOP counts efficiently. In addition, the experiment results in Section 5.1 show that the pruned ConvNets under this assumption are competitive in terms of performance with the pruned ConvNets obtained without this assumption. We state the subset assumption more formally below.

Assumption 1 (Subset Assumption) *For an optimal pruned ConvNet with FLOP count f , let $\mathcal{F}(f)_l$ be the filter count for layer l . The subset assumption states that $\mathcal{F}(f)_l \leq \mathcal{F}(f')_l \forall l$ if $f \leq f'$.*

Another challenge for learning a global ranking is the hardness of the problem. Obtaining an optimal global ranking can be expensive, *i.e.*, it requires $O(K \times K!)$ rounds of network fine-tuning, where K is the number of filters. Thus, to make it tractable, we assume the filter norm is able to rank filters locally (intra-layer-wise) but not globally (inter-layer-wise).

Assumption 2 (Norm Assumption) *ℓ_2 norm can be used to compare the importance of a filter within each layer, but not across layers.*

We note that the *norm assumption* is adopted and empirically verified by prior art [32, 62, 20]. For filter norms to be compared across layers, we propose to learn layer-wise affine transformations over filter norms. Specifically, the importance of filter i is defined as follows:

$$I_i = \alpha_{l(i)} \|\Theta_i\|_2^2 + \kappa_{l(i)}, \quad (1)$$

where $l(i)$ is the layer index for the i^{th} filter, $\|\cdot\|_2$ denotes ℓ_2 norms, Θ_i denotes the weights for the i^{th} filter, and $\alpha \in \mathbb{R}^L$, $\kappa \in \mathbb{R}^L$ are learnable parameters that represent layer-wise scale and shift values, and L denotes the number of layers. We will detail in Section 3.2 how α - κ pairs are learned so as to maximize overall accuracy.

Based on these learned affine transformations from Eq. (1) (*i.e.*, the α - κ pair), the *LeGR*-based pruning proceeds by ranking filters globally using I and prunes away bottom-ranked filters, *i.e.*, smaller in I , such that the FLOP count of interest is met, as shown in Figure 2. This process can be done efficiently without the need of training data (since the knowledge of pruning is encoded in the α - κ pair).

3.2. Learning Global Ranking

To learn α and κ , one can consider constructing a ranking with α and κ and then uniformly sampling ConvNets across different FLOP counts to evaluate the ranking. However, ConvNets obtained with different FLOP counts have drastically different validation accuracy, and one has to

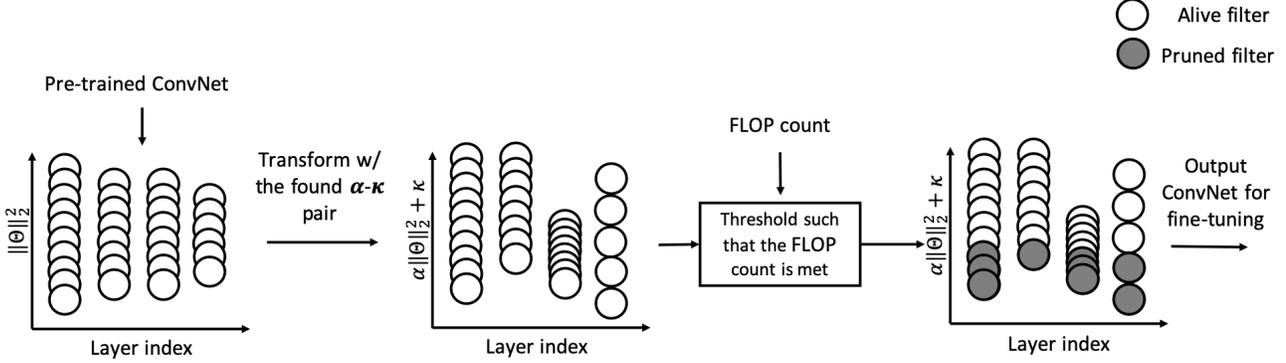


Figure 2: The flow of *LeGR-Pruning*. $\|\Theta\|_2^2$ represents the filter norm. Given the learned layer-wise affine transformations, *i.e.*, the α - κ pair, *LeGR-Pruning* returns filter masks that determine which filters are pruned. After *LeGR-Pruning*, the pruned network will be fine-tuned to obtain the final network.

know the Pareto curve³ of pruning to normalize the validation accuracy across ConvNets obtained with different FLOP counts. To address this difficulty, we propose to evaluate the validation accuracy of the ConvNet obtained from the lowest considered FLOP count as the objective for the ranking induced by the α - κ pair. Concretely, to learn α and κ , we treat *LeGR* as an optimization problem:

$$\arg \max_{\alpha, \kappa} Acc_{val}(\hat{\Theta}_l) \quad (2)$$

where

$$\hat{\Theta}_l = \text{LeGR-Pruning}(\alpha, \kappa, \hat{\zeta}_l). \quad (3)$$

LeGR-Pruning prunes away the bottom-ranked filters until the desired FLOP count is met as shown in Figure 2. $\hat{\zeta}_l$ denotes the lowest FLOP count considered. As we will discuss later in Section 5.1, we have also studied how $\hat{\zeta}$ affects the performance of the learned ranking, *i.e.*, how the learned ranking affects the accuracy of the pruned networks.

Specifically, to learn the α - κ pair, we rely on approaches from hyper-parameter optimization literature. While there are several options for the optimization algorithm, we adopt the regularized evolutionary algorithm (EA) proposed in [48] for its effectiveness in the neural architecture search space. The pseudo-code for our EA is outlined in Algorithm 1. We have also investigated policy gradients for solving for the α - κ pair, which is shown in Appendix B. We can equate each α - κ pair to a network architecture obtained by *LeGR-Pruning*. Once a pruned architecture is obtained, we fine-tune the resulting architecture by $\hat{\tau}$ gradient steps and use its accuracy on the validation set⁴ as the fitness (*i.e.*,

³A Pareto curve describes the optimal trade-off curve between two metrics of interest. Specifically, one cannot obtain improvement in one metric without degrading the other metric. The two metrics we considered in this work are accuracy and FLOP count.

⁴We split 10% of the original training set to be used as validation set.

Algorithm 1 Learning α, κ with regularized EA

Input: model Θ , lowest constraint $\hat{\zeta}_l$, random walk size σ , total search iterations E , sample size S , mutation ratio u , population size P , fine-tune iterations $\hat{\tau}$

Output: α, κ

Initialize *Pool* to a size P queue

for $e = 1$ **to** E **do**

$\alpha = 1, \kappa = 0$

if *Pool* has S samples **then**

$V = \text{Pool.sample}(S)$

$\alpha, \kappa = \text{argmaxFitness}(V)$

end if

$\text{Layer} = \text{Sample } u\% \text{ layers to mutate}$

for $l \in \text{Layer}$ **do**

$\text{std}_l = \text{computeStd}([M_i \forall i \in l])$

$\alpha_l = \alpha_l \times \hat{\alpha}_l$, where $\hat{\alpha}_l \sim e^{\mathcal{N}(0, \sigma^2)}$

$\kappa_l = \kappa_l + \hat{\kappa}_l$, where $\hat{\kappa}_l \sim \mathcal{N}(0, \text{std}_l)$

end for

$\hat{\Theta}_l = \text{LeGR-Pruning-and-fine-tuning}(\alpha, \kappa, \hat{\zeta}_l, \hat{\tau}, \Theta)$

$\text{Fitness} = Acc_{val}(\hat{\Theta}_l)$

$\text{Pool.replaceOldestWith}(\alpha, \kappa, \text{Fitness})$

end for

validation accuracy) for the corresponding α - κ pair. We note that we use $\hat{\tau}$ to approximate τ (fully fine-tuned steps) and we empirically find that $\hat{\tau} = 200$ gradient updates work well under the pruning settings across the datasets and networks we study. More concretely, we first generate a pool of candidates (α and κ values) and record the fitness for each candidate, and then repeat the following steps: (i) sample a subset from the candidates, (ii) identify the fittest candidate, (iii) generate a new candidate by mutating the fittest candidate and measure its fitness accordingly, and (iv) replace the oldest candidate in the pool with the generated one. To mu-

tate the fittest candidate, we randomly select a subset of the layers $Layer$ and conduct one step of random-walk from their current values, *i.e.*, $\alpha_l, \kappa_l \forall l \in Layer$.

We note that our layer-wise affine transformation formulation (Eq. 1) can be interpreted from an optimization perspective. That is, one can upper-bound the loss difference between a pre-trained ConvNet and its pruned-and-fine-tuned counterpart by assuming Lipschitz continuity on the loss function, as detailed in Appendix A.

4. Evaluations

4.1. Datasets and Training Setting

Our work is evaluated on various image classification benchmarks including CIFAR-10/100 [31], ImageNet [50], and Birds-200 [58]. CIFAR-10/100 consists of 50k training images and 10k testing images with a total of 10/100 classes to be classified. ImageNet is a large scale image classification dataset that includes 1.2 million training images and 50k testing images with 1k classes to be classified. Also, we benchmark the proposed algorithm in a transfer learning setting since in practice, we want a small and fast model on some target datasets. Specifically, we use the Birds-200 dataset that consists of 6k training images and 5.7k testing images covering 200 bird species.

For Bird-200, we use 10% of the training data as the validation set used for early stopping and to avoid over-fitting. The training scheme for CIFAR-10/100 follows [19], which uses stochastic gradient descent with nesterov [44], weight decay $5e^{-4}$, batch size 128, $1e^{-1}$ initial learning rate with decrease by $5\times$ at epochs 60, 120, and 160, and train for 200 epochs in total. For control experiments with CIFAR-100 and Bird-200, the fine-tuning after pruning is done as follows: we keep all training hyper-parameters the same but change the initial learning rate to $1e^{-2}$ and train for 60 epochs (*i.e.*, $\tau \approx 21k$). We drop the learning rate by $10\times$ at 30%, 60%, and 80% of the total epochs, *i.e.*, epochs 18, 36, and 48. To compare numbers with prior art on CIFAR-10 and ImageNet, we follow the number of iterations in [72]. Specifically, for CIFAR-10 we fine-tuned for 400 epochs with initial learning rate $1e^{-2}$, drop by $5\times$ at epochs 120, 240, and 320. For ImageNet, we use pre-trained models and we fine-tuned the pruned models for 60 epochs with initial learning rate $1e^{-2}$, drop by $10\times$ at epochs 30 and 45.

For the hyper-parameters of *LeGR*, we select $\hat{\tau} = 200$, *i.e.*, fine-tune for 200 gradient steps before measuring the validation accuracy when searching for the α - κ pair. We note that we do the same for *AMC* [20] for a fair comparison. Moreover, we set the number of architectures explored to be the same with *AMC*, *i.e.*, 400. We set mutation rate $u = 10$ and the hyper-parameter of the regularized evolutionary algorithm by following prior art [48]. In the following experiments, we use the smallest ζ considered as $\hat{\zeta}_l$

to search for the learnable variables α and κ . The found α - κ pair is used to obtain the pruned networks at various FLOP counts. For example, for ResNet-56 with CIFAR-100 (Figure 3a), we use $\hat{\zeta}_l = 20\%$ to obtain the α - κ pair and use the same α - κ pair to obtain the seven networks ($\zeta = 20\%, \dots, 80\%$) with the flow described in Figure 2. The ablation of $\hat{\zeta}_l$ and $\hat{\tau}$ are detailed in Sec. 5.2.

We prune filters across all the convolutional layers. We group dependent channels by summing up their importance measure and prune them jointly. The importance measure refers to the measure after learned affine transformations. Specifically, we group a channel in depth-wise convolution with its corresponding channel in the preceding layer. We also group channels that are summed together through residual connections.

4.2. CIFAR-100 Results

In this section, we consider ResNet-56 and MobileNetV2 and we compare *LeGR* mainly with four filter pruning methods, *i.e.*, *MorphNet* [16], *AMC* [20], *FisherPruning* [57], and a baseline that prunes filters uniformly across layers. Specifically, the baselines are determined such that one dominant approach is selected from different groups of prior art. We select one approach [16] from pruning-while-learning approaches, one approach [20] from pruning-by-searching methods, one approach [57] from continuous pruning methods, and a baseline extending magnitude-based pruning to various FLOP counts. We note that *FisherPruning* is a continuous pruning method where we use 0.0025 learning rate and perform 500 gradient steps after each filter pruned following [57].

As shown in Figure 3a, we first observe that *FisherPruning* does not work as well as other methods and we hypothesize the reason for it is that the small fixed learning rate in the fine-tuning phase makes it hard for the optimizer to get out of local optima. Additionally, we find that *FisherPruning* prunes away almost all the filters for some layers. On the other hand, we find that all other approaches outperform the uniform baseline in a high-FLOP-count regime. However, both *AMC* and *MorphNet* have higher variances when pruned more aggressively. In both cases, *LeGR* outperforms prior art, especially in the low-FLOP-count regime.

More importantly, our proposed method aims to alleviate the cost of pruning when the goal is to explore the trade-off curve between accuracy and inference latency. From this perspective, our approach outperforms prior art by a significant margin. More specifically, we measure the average time of each algorithm to obtain the seven pruned ResNet-56 across the FLOP counts in Figure 3a using our hardware (*i.e.*, NVIDIA GTX 1080 Ti). Figure 3b shows the efficiency of *AMC*, *MorphNet*, *FisherPruning*, and the proposed *LeGR*. The cost can be broken down into two parts: (1) pruning: the time it takes to search for a network that has

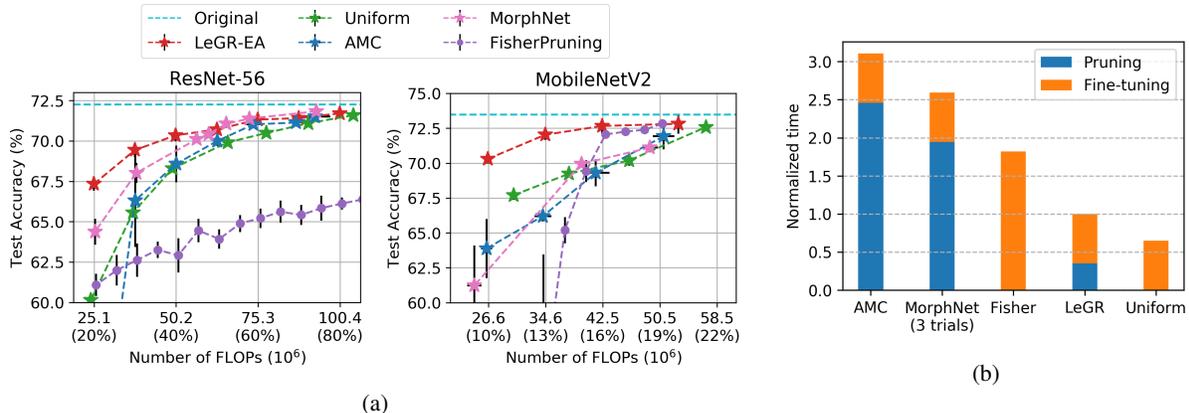


Figure 3: (a) The trade-off curve of pruning ResNet-56 and MobileNetV2 on CIFAR-100 using various methods. We average across three trials and plot the mean and standard deviation. (b) Training cost for seven ConvNets across FLOP counts using various methods targeting ResNet-56 on CIFAR-100. We report the average cost considering seven FLOP counts, *i.e.*, 20% to 80% FLOP count in a step of 10% on NVIDIA GTX 1080 Ti. The cost is normalized to the cost of LeGR.

some pre-defined FLOP count and (2) fine-tuning: the time it takes for fine-tuning the weights of a pruned network. For *MorphNet*, we consider three trials for each FLOP count to find an appropriate hyper-parameter λ to meet the FLOP count of interest. The numbers are normalized to the cost of *LeGR*. In terms of pruning time, *LeGR* is $7\times$ and $5\times$ faster than *AMC* and *MorphNet*, respectively. The efficiency comes from the fact that *LeGR* only searches the α - κ pair once and re-uses it across FLOP counts. In contrast, both *AMC* and *MorphNet* have to search for networks for every FLOP count considered. *FisherPruning* always prune one filter at a time, and therefore the lowest FLOP count level considered determines the pruning time, regardless of how many FLOP count levels we are interested in.

4.3. Comparison with Prior Art

Although the goal of this work is to develop a model compression method that produces a set of ConvNets across different FLOP counts, we also compare our method with prior art that focuses on generating a ConvNet for a specified FLOP count.

CIFAR-10 In Table 1, we compare *LeGR* with prior art that reports results on CIFAR-10. First, for ResNet-56, we find that *LeGR* outperforms most of the prior art in both FLOP count and accuracy dimensions and performs similarly to [19, 72]. For VGG-13, *LeGR* achieves significantly better results compared to prior art.

ImageNet Results For ImageNet, we prune ResNet-50 and MobileNetV2 with *LeGR* to compare with prior art. For *LeGR*, we learn the ranking using 47% FLOP count for ResNet-50 and 50% FLOP count for MobileNetV2, and use the learned ranking to obtain ConvNets for other FLOP

Table 1: Comparison with prior art on CIFAR-10. We group methods into sections according to different FLOP counts. Values for our approaches are averaged across three trials and we report the mean and standard deviation. We use boldface to denote the best numbers and use * to denote our implementation. The accuracy is represented in the format of *pre-trained* \mapsto *pruned-and-fine-tuned*.

NETWORK	METHOD	Acc. (%)	MFLOP COUNT
RESNET-56	PF [32]	93.0 \rightarrow 93.0	90.9 (72%)
	TAYLOR [43]*	93.9 \rightarrow 93.2	90.8 (72%)
	LEGR	93.9 \rightarrow 94.1\pm0.0	87.8 (70%)
	DCP-ADAPT [72]	93.8 \rightarrow 93.8	66.3 (53%)
	CP [23]	92.8 \rightarrow 91.8	62.7 (50%)
	AMC [20]	92.8 \rightarrow 91.9	62.7 (50%)
	DCP [72]	93.8 \rightarrow 93.5	62.7 (50%)
	SFP [19]	93.6 \pm 0.6 \rightarrow 93.4\pm0.3	59.4 (47%)
	LEGR	93.9 \rightarrow 93.7\pm0.2	58.9 (47%)
VGG-13	BC-GNJ [38]	91.9 \rightarrow 91.4	141.5 (45%)
	BC-GHS [38]	91.9 \rightarrow 91	121.9 (39%)
	VIBNET [8]	91.9 \rightarrow 91.5	70.6 (22%)
	LEGR	91.9 \rightarrow 92.4\pm0.2	70.3 (22%)

counts of interest. We have compared to 17 prior methods that report pruning performance for ResNet-50 and/or MobileNetV2 on the ImageNet dataset. While our focus is on the fast exploration of the speed and accuracy trade-off curve for filter pruning, our proposed method is better or comparable compared to the state-of-the-art methods as shown in Figure 4. The detailed numerical results are in Appendix C. We would like to emphasize that to obtain a pruned-ConvNet with prior methods, one has to run the pruning algorithm for every FLOP count considered. In contrast, our proposed method learns the ranking once and uses it to obtain ConvNets across different FLOP counts.

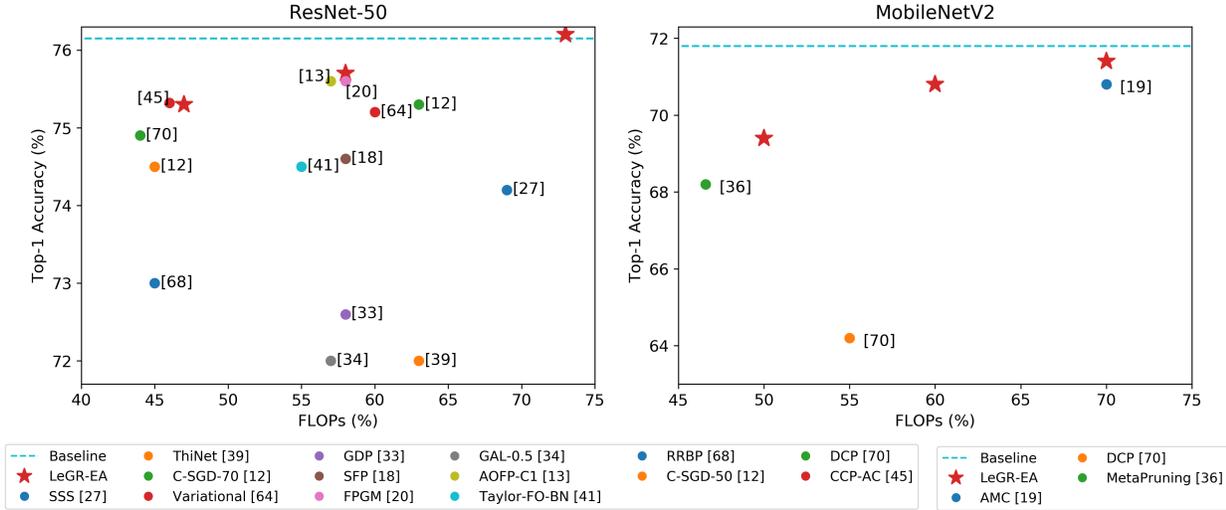


Figure 4: Results for ImageNet. LeGR is better or comparable compared to prior methods. Furthermore, its goal is to output a set of ConvNets instead of one ConvNet. The detailed numerical results are in Appendix C.

4.4. Transfer Learning: Bird-200

We analyze how *LeGR* performs in a transfer learning setting where we have a model pre-trained on a large dataset, *i.e.*, ImageNet, and we want to transfer its knowledge to adapt to a smaller dataset, *i.e.*, Bird-200. We prune the fine-tuned network on the target dataset directly following the practice in prior art [68, 40]. We first obtain fine-tuned MobileNetV2 and ResNet-50 on the Bird-200 dataset with top-1 accuracy 80.2% and 79.5%, respectively. These are comparable to the reported values in prior art [33, 41]. As shown in Figure 5, we find that *LeGR* outperforms *Uniform* and *AMC*, which is consistent with previous analyses in Section 4.2.

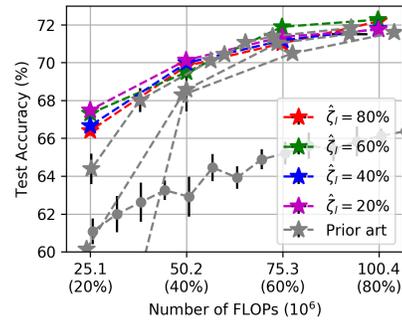


Figure 6: Robustness to the hyper-parameter $\hat{\zeta}_l$. Prior art is plotted as a reference (c.f. Figure 3a).

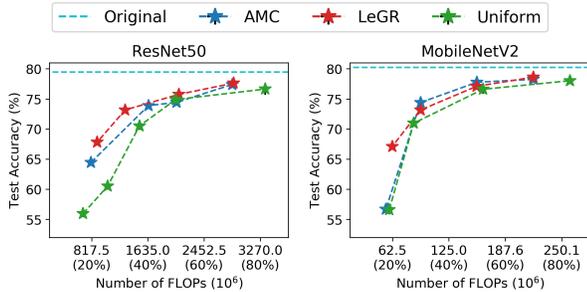


Figure 5: Results for Bird-200.

5. Ablation Study

5.1. Ranking Performance and $\hat{\zeta}_l$

To learn the global ranking with *LeGR* without knowing the Pareto curve in advance, we use the minimum consid-

ered FLOP count ($\hat{\zeta}_l$) during learning to evaluate the performance of a ranking. We are interested in understanding how this design choice affects the performance of *LeGR*. Specifically, we try *LeGR* targeting ResNet-56 for CIFAR-100 with $\hat{\zeta}_l \in \{20\%, 40\%, 60\%, 80\%\}$. As shown in Figure 6, we first observe that rankings learned using different FLOP counts have similar performances, which empirically supports Assumption 1. More concretely, consider the network pruned to 40% FLOP count by using the ranking learned at 40% FLOP count. This case does not take advantage of the subset assumption because the entire learning process for learning $\alpha\text{-}\kappa$ is done only by looking at the performance of the 40% FLOP count network. On the other hand, rankings learned using other FLOP counts but employed to obtain pruned-networks at 40% FLOP count have exploited the subset assumption (*e.g.*, the ranking learned for 80% FLOP count can produce a competitive network

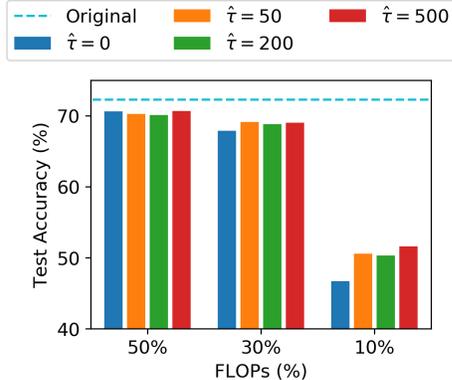


Figure 7: Pruning ResNet-56 for CIFAR-100 with *LeGR* by learning α and κ using different $\hat{\tau}$ and FLOP count constraints.

for 40% FLOP count). We find that *LeGR* with or without employing Assumption 1 results in similar performance for the pruned networks.

5.2. Fine-tuned Iterations

Since we use $\hat{\tau}$ to approximate τ when learning the α - κ pair, it is expected that the closer $\hat{\tau}$ to τ , the better the α - κ pair *LeGR* can find. We use *LeGR* to prune ResNet-56 for CIFAR-100 and learn α - κ at three FLOP counts $\hat{\zeta}_l \in \{10\%, 30\%, 50\%\}$. We consider ζ to be exactly $\hat{\zeta}_l$ in this case. For $\hat{\tau}$, we experiment with $\{0, 50, 200, 500\}$. We note that once the α - κ pair is learned, we use *LeGR-Pruning* to obtain the pruned ConvNet, fine-tune it for τ steps, and plot the resulting test accuracy. In this experiment, τ is set to 21120 gradient steps (60 epochs). As shown in Figure 7, the results align with our intuition in that there are diminishing returns in increasing $\hat{\tau}$. We observe that $\hat{\tau}$ affects the accuracy of the pruned ConvNets more when learning the ranking at a lower FLOP count level, which means in low-FLOP-count regimes, the validation accuracy after fine-tuning a few steps might not be representative. This makes sense since when pruning away a lot of filters, the network can be thought of as moving far away from the local optimal, where the gradient steps early in the fine-tuning phase are noisy. Thus, more gradient steps are needed before considering the accuracy to be representative of the fully-fine-tuned accuracy.

5.3. FLOP count and Runtime

We demonstrate the effectiveness of filter pruning in wall-clock time speedup using ResNet-50 and MobileNetV2 on PyTorch 0.4 using two types of CPUs. Specifically, we consider both a desktop level CPU, *i.e.*, Intel i7, and an embedded CPU, *i.e.*, ARM A57, and use *LeGR* as the pruning methodology. The input is a single RGB image

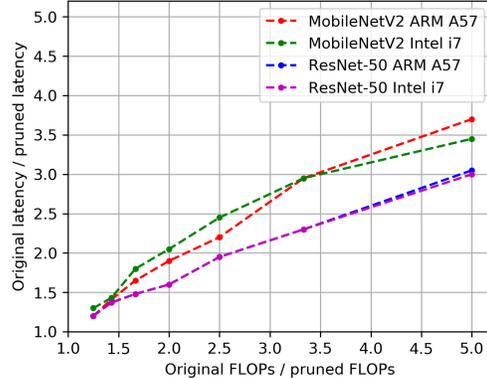


Figure 8: Latency reduction vs. FLOP count reduction. FLOP count reduction is indicative for latency reduction.

of size 224x224 and the program (Python with PyTorch) is run using a single thread. As shown in Figure 8, filter pruning can produce near-linear acceleration (with a slope of approximately 0.6) without specialized software or hardware support.

6. Conclusion

To alleviate the bottleneck of using model compression in optimizing the ConvNets in a large system, we propose *LeGR*, a novel formulation for practitioners to explore the accuracy-vs-speed trade-off efficiently via filter pruning. More specifically, we propose to learn layer-wise affine transformations over filter norms to construct a global ranking of filters. This formulation addresses the limitation that filter norms cannot be compared across layers in a learnable fashion and provides an efficient way for practitioners to obtain ConvNet architectures with different FLOP counts. Additionally, we provide a theoretical interpretation of the proposed affine transformation formulation. We conduct extensive empirical analyses using ResNet and MobileNetV2 on datasets including CIFAR, Bird-200, and ImageNet and show that *LeGR* has less training cost to generate the pruned ConvNets across different FLOP counts compared to prior art while achieving comparable performance to state-of-the-art pruning methods.

Acknowledgement

This research was supported in part by NSF CCF Grant No. 1815899, NSF CSR Grant No. 1815780, and NSF ACI Grant No. 1445606 at the Pittsburgh Supercomputing Center (PSC).

References

- [1] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhuf, Ian Reid, Stephen Gould, and

- Anton van den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3674–3683, 2018. 1
- [2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015. 1
- [3] Behzad Boroujerdian, Hasan Genc, Srivatsan Krishnan, Wenzhi Cui, Aleksandra Faust, and Vijay Reddi. Mavbench: Micro aerial vehicle benchmarking. In *2018 51st Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pages 894–907. IEEE, 2018. 1
- [4] Han Cai, Ligeng Zhu, and Song Han. Proxylessnas: Direct neural architecture search on target task and hardware. *arXiv preprint arXiv:1812.00332*, 2018. 2
- [5] Ting-Wu Chin, Ruizhou Ding, and Diana Marculescu. Adascale: Towards real-time video object detection using adaptive scaling. *arXiv preprint arXiv:1902.02910*, 2019. 2
- [6] Chung-Cheng Chiu, Tara N Sainath, Yonghui Wu, Rohit Prabhavalkar, Patrick Nguyen, Zhifeng Chen, Anjali Kannan, Ron J Weiss, Kanishka Rao, Ekaterina Gonina, et al. State-of-the-art speech recognition with sequence-to-sequence models. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4774–4778. IEEE, 2018. 1
- [7] Jungwook Choi, Pierce I-Jen Chuang, Zhuo Wang, Swagath Venkataramani, Vijayalakshmi Srinivasan, and Kailash Gopalakrishnan. Bridging the accuracy gap for 2-bit quantized neural networks (qnn). *arXiv preprint arXiv:1807.06964*, 2018. 2
- [8] Bin Dai, Chen Zhu, and David Wipf. Compressing neural networks using the variational information bottleneck. *arXiv preprint arXiv:1802.10399*, 2018. 1, 6
- [9] Xiaoliang Dai, Peizhao Zhang, Bichen Wu, Hongxu Yin, Fei Sun, Yanghan Wang, Marat Dukhan, Yunqing Hu, Yiming Wu, Yangqing Jia, et al. Chamnet: Towards efficient network design through platform-aware model adaptation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [10] Zihang Dai, Zhilin Yang, Yiming Yang, William W Cohen, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019. 1
- [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 1
- [12] Ruizhou Ding, Ting-Wu Chin, Zeye Liu, and Diana Marculescu. Regularizing activation distribution for training binarized deep networks. 2019. 2
- [13] Xiaohan Ding, Guiguang Ding, Yuchen Guo, and Jungong Han. Centripetal sgd for pruning very deep convolutional networks with complicated structure. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4943–4953, 2019. 13
- [14] Xiaohan Ding, Guiguang Ding, Yuchen Guo, Jungong Han, and Chenggang Yan. Approximated oracle filter pruning for destructive cnn width optimization. *arXiv preprint arXiv:1905.04748*, 2019. 13
- [15] Jin-Dong Dong, An-Chieh Cheng, Da-Cheng Juan, Wei Wei, and Min Sun. Dpp-net: Device-aware progressive search for pareto-optimal neural architectures. *arXiv preprint arXiv:1806.08198*, 2018. 2
- [16] Ariel Gordon, Elad Eban, Ofir Nachum, Bo Chen, Hao Wu, Tien-Ju Yang, and Edward Choi. Morphnet: Fast & simple resource-constrained structure learning of deep networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1, 2, 5
- [17] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 1
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1
- [19] Yang He, Guoliang Kang, Xuanyi Dong, Yanwei Fu, and Yi Yang. Soft filter pruning for accelerating deep convolutional neural networks. In *IJCAI*, pages 2234–2240, 2018. 2, 5, 6, 13
- [20] Yihui He, Ji Lin, Zhijian Liu, Hanrui Wang, Li-Jia Li, and Song Han. Amc: Automl for model compression and acceleration on mobile devices. *arXiv preprint arXiv:1802.03494*, 2018. 1, 2, 3, 5, 6, 12, 13
- [21] Yang He, Ping Liu, Ziwei Wang, Zhilan Hu, and Yi Yang. Filter pruning via geometric median for deep convolutional neural networks acceleration. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4340–4349, 2019. 13
- [22] Yihui He, Xianggen Liu, Huasong Zhong, and Yuchun Ma. Addressnet: Shift-based primitives for efficient convolutional neural networks. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1213–1222. IEEE, 2019. 2
- [23] Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks. In *International Conference on Computer Vision (ICCV)*, volume 2, 2017. 2, 6, 13
- [24] Lu Hou and James T. Kwok. Loss-aware weight quantization of deep networks. In *International Conference on Learning Representations*, 2018. 2
- [25] 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 preprint arXiv:1704.04861*, 2017. 2
- [26] Gao Huang, Shichen Liu, Laurens van der Maaten, and Kilian Q Weinberger. Condensenet: An efficient densenet using learned group convolutions. *group*, 3(12):11, 2017. 2
- [27] Zehao Huang and Naiyan Wang. Data-driven sparse structure selection for deep neural networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 304–320, 2018. 2

- [28] Zehao Huang and Naiyan Wang. Data-driven sparse structure selection for deep neural networks. In *The European Conference on Computer Vision (ECCV)*, September 2018. 13
- [29] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2
- [30] Sangil Jung, Changyong Son, Seohyung Lee, Jinwoo Son, Jae-Joon Han, Youngjun Kwak, Sung Ju Hwang, and Changkyu Choi. Learning to quantize deep networks by optimizing quantization intervals with task loss. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 2
- [31] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009. 5
- [32] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. *International Conference on Learning Representations (ICLR)*, 2017. 1, 2, 3, 6
- [33] Xingjian Li, Haoyi Xiong, Hanchao Wang, Yuxuan Rao, Liping Liu, and Jun Huan. DELTA: DEEP LEARNING TRANSFER USING FEATURE MAP WITH ATTENTION FOR CONVOLUTIONAL NETWORKS. In *International Conference on Learning Representations*, 2019. 7
- [34] Shaohui Lin, Rongrong Ji, Yuchao Li, Yongjian Wu, Feiyue Huang, and Baochang Zhang. Accelerating convolutional networks via global & dynamic filter pruning. In *IJCAI*, pages 2425–2432, 2018. 13
- [35] Shaohui Lin, Rongrong Ji, Chenqian Yan, Baochang Zhang, Liujuan Cao, Qixiang Ye, Feiyue Huang, and David Doermann. Towards optimal structured cnn pruning via generative adversarial learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2790–2799, 2019. 13
- [36] Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang. Learning efficient convolutional networks through network slimming. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 2755–2763. IEEE, 2017. 1, 2
- [37] Zechun Liu, Haoyuan Mu, Xiangyu Zhang, Zichao Guo, Xin Yang, Tim Kwang-Ting Cheng, and Jian Sun. Metapruning: Meta learning for automatic neural network channel pruning. In *Proceedings of the IEEE International Conference on Computer Vision*, 2019. 13
- [38] Christos Louizos, Karen Ullrich, and Max Welling. Bayesian compression for deep learning. In *Advances in Neural Information Processing Systems*, pages 3288–3298, 2017. 2, 6
- [39] Christos Louizos, Max Welling, and Diederik P. Kingma. Learning sparse neural networks through l_0 regularization. In *International Conference on Learning Representations*, 2018. 2
- [40] Jian-Hao Luo, Jianxin Wu, and Weiyao Lin. Thinet: A filter level pruning method for deep neural network compression. *arXiv preprint arXiv:1707.06342*, 2017. 2, 7, 13
- [41] Arun Mallya and Svetlana Lazebnik. Packnet: Adding multiple tasks to a single network by iterative pruning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 7
- [42] Pavlo Molchanov, Arun Mallya, Stephen Tyree, Iuri Frosio, and Jan Kautz. Importance estimation for neural network pruning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 11264–11272, 2019. 13
- [43] Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning convolutional neural networks for resource efficient inference. *International Conference on Learning Representation (ICLR)*, 2017. 3, 6
- [44] Yurii E Nesterov. A method for solving the convex programming problem with convergence rate $o(1/k^2)$. In *Dokl. Akad. Nauk SSSR*, volume 269, pages 543–547, 1983. 5
- [45] Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk, and Quoc V Le. SpecAugment: A simple data augmentation method for automatic speech recognition. *arXiv preprint arXiv:1904.08779*, 2019. 1
- [46] Hanyu Peng, Jiayang Wu, Shifeng Chen, and Junzhou Huang. Collaborative channel pruning for deep networks. In *International Conference on Machine Learning*, pages 5113–5122, 2019. 1, 13
- [47] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. In *European Conference on Computer Vision*, pages 525–542. Springer, 2016. 2
- [48] Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V Le. Regularized evolution for image classifier architecture search. *arXiv preprint arXiv:1802.01548*, 2018. 4, 5
- [49] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. 1
- [50] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. 5
- [51] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A platform for embodied ai research. *arXiv preprint arXiv:1904.01201*, 2019. 1
- [52] Dimitrios Stamoulis, Ting-Wu Rudy Chin, Anand Krishnan Prakash, Haocheng Fang, Sribhuvan Sajja, Mitchell Bogner, and Diana Marculescu. Designing adaptive neural networks for energy-constrained image classification. In *Proceedings of the International Conference on Computer-Aided Design*, page 23. ACM, 2018. 2
- [53] Dimitrios Stamoulis, Ruizhou Ding, Di Wang, Dimitrios Lymberopoulos, Bodhi Priyantha, Jie Liu, and Diana Marculescu. Single-path nas: Designing hardware-efficient convnets in less than 4 hours. *arXiv preprint arXiv:1904.02877*, 2019. 2

- [54] Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, and Quoc V Le. Mnasnet: Platform-aware neural architecture search for mobile. *arXiv preprint arXiv:1807.11626*, 2018. [2](#)
- [55] Mingxing Tan and Quoc Le. EfficientNet: Rethinking model scaling for convolutional neural networks. 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 6105–6114, Long Beach, California, USA, 09–15 Jun 2019. PMLR. [2](#)
- [56] Makarand Tapaswi, Yukun Zhu, Rainer Stiefelhagen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. Movieqa: Understanding stories in movies through question-answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4631–4640, 2016. [1](#)
- [57] Lucas Theis, Iryna Korshunova, Alykhan Tejani, and Ferenc Huszár. Faster gaze prediction with dense networks and fisher pruning. *arXiv preprint arXiv:1801.05787*, 2018. [3](#), [5](#)
- [58] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011. [5](#)
- [59] Huan Wang, Qiming Zhang, Yuehai Wang, and Haoji Hu. Structured probabilistic pruning for convolutional neural network acceleration. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2018. [13](#)
- [60] Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity in deep neural networks. In *Advances in Neural Information Processing Systems*, pages 2074–2082, 2016. [1](#), [2](#)
- [61] Bichen Wu, Alvin Wan, Xiangyu Yue, Peter Jin, Sicheng Zhao, Noah Golmant, Amir Gholaminejad, Joseph Gonzalez, and Kurt Keutzer. Shift: A zero flop, zero parameter alternative to spatial convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9127–9135, 2018. [2](#)
- [62] Tien-Ju Yang, Andrew Howard, Bo Chen, Xiao Zhang, Alec Go, Vivienne Sze, and Hartwig Adam. Netadapt: Platform-aware neural network adaptation for mobile applications. *arXiv preprint arXiv:1804.03230*, 2018. [3](#)
- [63] Jianbo Ye, Xin Lu, Zhe Lin, and James Z Wang. Rethinking the smaller-norm-less-informative assumption in channel pruning of convolution layers. *International Conference on Learning Representation (ICLR)*, 2018. [1](#)
- [64] Ruichi Yu, Ang Li, Chun-Fu Chen, Jui-Hsin Lai, Vlad I. Morariu, Xintong Han, Mingfei Gao, Ching-Yung Lin, and Larry S. Davis. Nisp: Pruning networks using neuron importance score propagation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. [13](#)
- [65] Xin Yuan, Liangliang Ren, Jiwen Lu, and Jie Zhou. Enhanced bayesian compression via deep reinforcement learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. [2](#)
- [66] Chenglong Zhao, Bingbing Ni, Jian Zhang, Qiwei Zhao, Wenjun Zhang, and Qi Tian. Variational convolutional neural network pruning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2780–2789, 2019. [13](#)
- [67] Ritchie Zhao, Yuwei Hu, Jordan Dotzel, Christopher De Sa, and Zhiru Zhang. Building efficient deep neural networks with unitary group convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 11303–11312, 2019. [2](#)
- [68] Yang Zhong, Vladimir Li, Ryuzo Okada, and Atsuto Maki. Target aware network adaptation for efficient representation learning. *arXiv preprint arXiv:1810.01104*, 2018. [7](#)
- [69] Yanqi Zhou, Siavash Ebrahimi, Sercan Ö Arik, Haonan Yu, Hairong Liu, and Greg Diamos. Resource-efficient neural architect. *arXiv preprint arXiv:1806.07912*, 2018. [2](#)
- [70] Yuefu Zhou, Ya Zhang, Yanfeng Wang, and Qi Tian. Accelerate cnn via recursive bayesian pruning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3306–3315, 2019. [1](#), [2](#), [13](#)
- [71] Chenzhuo Zhu, Song Han, Huizi Mao, and William J Dally. Trained ternary quantization. In *International Conference on Learning Representations*, 2017. [2](#)
- [72] Zhuangwei Zhuang, Mingkui Tan, Bohan Zhuang, Jing Liu, Yong Guo, Qingyao Wu, Junzhou Huang, and Jinhui Zhu. Discrimination-aware channel pruning for deep neural networks. In *Advances in Neural Information Processing Systems*, pages 883–894, 2018. [5](#), [6](#), [13](#)

A. Optimization Interpretation of *LeGR*

LeGR can be interpreted as minimizing a surrogate of a derived upper bound for the loss difference between (1) the pruned-and-fine-tuned CNN and (2) the pre-trained CNN. Concretely, we would like to solve for the filter masking binary variables $\mathbf{z} \in \{0, 1\}^K$, with K being the number of filters. If a filter k is pruned, the corresponding mask will be zero ($z_k = 0$), otherwise it will be one ($z_k = 1$). Thus, we have the following optimization problem:

$$\begin{aligned} \min_{\mathbf{z}} \mathbb{L}(\Theta \odot \mathbf{z} - \eta \sum_{j=1}^{\tau} \Delta \mathbf{w}^{(j)} \odot \mathbf{z}) - \mathbb{L}(\Theta) \\ \text{s.t. } C(\mathbf{z}) \leq \zeta, \end{aligned} \quad (4)$$

where Θ denotes all the filters of the CNN, $\mathbb{L}(\Theta) = \frac{1}{|D|} \sum_{(x,y) \in D} L(f(x|\Theta), y)$ denotes the loss function of filters where x and y are the input and label, respectively. D denotes the training data, f is the CNN model and L is the loss function for prediction (e.g., cross entropy loss). η denotes the learning rate, τ denotes the number of gradient steps, $\Delta \mathbf{w}^{(j)}$ denotes the gradient with respect to the filter weights computed at step j , and \odot denotes element-wise multiplication. On the constraint side, $C(\cdot)$ is the modeling function for FLOP count and ζ is the desired FLOP count constraint. By fine-tuning, we mean updating the filter weights with stochastic gradient descent (SGD) for τ steps.

Let us assume the loss function \mathbb{L} is Ω_l -Lipschitz continuous for the l -th layer of the CNN, then the following holds:

$$\begin{aligned} & \mathbb{L}(\Theta \odot \mathbf{z} - \eta \sum_{j=1}^{\tau} \Delta \mathbf{w}^{(j)} \odot \mathbf{z}) - \mathbb{L}(\Theta) \\ & \leq \mathbb{L}(\Theta \odot \mathbf{z}) + \sum_{i=1}^K \Omega_{l(i)} \eta \left\| \sum_{j=1}^{\tau} \Delta \mathbf{w}_i^{(j)} \odot \mathbf{z}_i \right\| - \mathbb{L}(\Theta) \\ & \leq \sum_{i=1}^K \Omega_{l(i)} \|\Theta_i\| \mathbf{h}_i + \sum_{i=1}^K \Omega_{l(i)}^2 \eta \tau \mathbf{z}_i \\ & = \sum_{i=1}^K (\Omega_{l(i)} \|\Theta_i\| - \Omega_{l(i)}^2 \eta \tau) \mathbf{h}_i + \Omega_{l(i)}^2 \eta \tau, \end{aligned} \quad (5)$$

where $l(i)$ is the layer index for the i -th filter, $\mathbf{h} = \mathbf{1} - \mathbf{z}$, and $\|\cdot\|$ denotes ℓ_2 norms.

On the constraint side of equation (4), let $R_{l(i)}$ be the FLOP count of layer $l(i)$ where filter i resides. Analytically, the FLOP count of a layer depends linearly on the number of filters in its preceding layer:

$$R_{l(i)} = u_{l(i)} \|\{\mathbf{z} : \mathbf{z}_j \forall j \in P(l(i))\}\|_0, \quad u_{l(i)} \geq 0, \quad (6)$$

where $P(l(i))$ returns a set of filter indices for the layer that precedes layer $l(i)$ and $u_{l(i)}$ is a layer-dependent positive constant. Let $\hat{R}_{l(i)}$ denote the FLOP count for layer $l(i)$ for the pre-trained network ($\mathbf{z} = \mathbf{1}$), one can see from equation (6) that $R_{l(i)} \leq \hat{R}_{l(i)} \forall i, \mathbf{z}$. Thus, the following holds:

$$C(\mathbf{1} - \mathbf{h}) = \sum_i^K R_{l(i)} (1 - \mathbf{h}_i) \leq \sum_i^K \hat{R}_{l(i)} (1 - \mathbf{h}_i). \quad (7)$$

Based on equations (5) and (7), instead of minimizing equation (4), we minimize its upper bound in a Lagrangian form. That is,

$$\min_{\mathbf{h}} \sum_{i=1}^K (\alpha_{l(i)} \|\Theta_i\| + \kappa_{l(i)}) \mathbf{h}_i, \quad (8)$$

where $\alpha_{l(i)} = \Omega_{l(i)}$ and $\kappa_{l(i)} = \eta \tau \Omega_{l(i)}^2 - \lambda \hat{R}_{l(i)}$. To guarantee the solution will satisfy the constraint, we rank all filters by their scores $s_i = \alpha_{l(i)} \|\Theta_i\| + \kappa_{l(i)} \forall i$ and threshold out the bottom ranked (small in scores) filters such that the constraint $C(\mathbf{1} - \mathbf{h}) \leq \zeta$ is satisfied and $\|\mathbf{h}\|_0$ is maximized. That is, *LeGR* can be viewed as learning to estimate α and κ by assuming that better estimates of α - κ produce a better solution for the original objective (4) by solving the surrogate of the upper bound (8).

B. LeGR-DDPG

We have also tried learning the layer-wise affine transformations with actor-critic policy gradient (DDPG), which is adopted in prior art [20]. We use DDPG in a sequential fashion that follows [20]. *LeGR* requires two continuous actions (i.e., α_l and κ_l) for layer l while *AMC* needs only one action (i.e., percentage). We conduct the comparison of pruning ResNet-56 to 50% of its original FLOP count targeting CIFAR-100 with $\hat{\tau} = 0$ and hyper-parameters following [20]. As shown in Fig. 9a, while both *LeGR* and *AMC* outperform random search (iterations before the vertical black-dotted line), *LeGR* converges faster to a better solution. Beyond comparing the progress of searching, we also compare the performance of the final pruned networks. As shown in Fig. 9b, searching layer-wise affine transformations is more efficient and effective compared to searching the layer-wise filter percentages. Comparing *LeGR* using the two policy improvement methods, we empirically find that DDPG incurs larger variance on the final network than evolutionary algorithm.

C. ImageNet Result Detail

The comparison of *LeGR* with prior art on ImageNet is detailed in Table 2.

Table 2: Summary of pruning on ImageNet. The sections are defined based on the FLOP count left. The accuracy is represented in the format of *pre-trained* \mapsto *pruned-and-fine-tuned*.

NETWORK	METHOD	TOP-1	TOP-1 DIFF	TOP-5	TOP-5 DIFF	FLOP COUNT (%)
RESNET-50	NISP [64]	- \rightarrow -	-0.2	- \rightarrow -	-	73
	LEGR	76.1 \rightarrow 76.2	+0.1	92.9 \rightarrow 93.0	+0.1	73
	SSS [28]	76.1 \rightarrow 74.2	-1.9	92.9 \rightarrow 91.9	-1.0	69
	THiNET [40]	72.9 \rightarrow 72.0	-0.9	91.1 \rightarrow 90.7	-0.4	63
	C-SGD-70 [13]	75.3 \rightarrow 75.3	+0.0	92.6 \rightarrow 92.5	-0.1	63
	VARIATIONAL [66]	75.1 \rightarrow 75.2	+0.1	92.8 \rightarrow 92.1	-0.7	60
	GDP [34]	75.1 \rightarrow 72.6	-2.5	92.3 \rightarrow 91.1	-1.2	58
	SFP [19]	76.2 \rightarrow 74.6	-1.6	92.9 \rightarrow 92.1	-0.8	58
	FPGM [21]	76.2 \rightarrow 75.6	-0.6	92.9 \rightarrow 92.6	-0.3	58
	LEGR	76.1 \rightarrow 75.7	-0.4	92.9 \rightarrow 92.7	-0.2	58
	GAL-0.5 [35]	76.2 \rightarrow 72.0	-4.2	92.9 \rightarrow 91.8	-1.1	57
	AOFP-C1 [14]	75.3 \rightarrow 75.6	+0.3	92.6 \rightarrow 92.7	+0.1	57
	NISP [64]	- \rightarrow -	-0.9	- \rightarrow -	-	56
	TAYLOR-FO-BN [42]	76.2 \rightarrow 74.5	-1.7	- \rightarrow -	-	55
	CP [23]	- \rightarrow -	-	92.2 \rightarrow 90.8	-1.4	50
	SPP [59]	- \rightarrow -	-	91.2 \rightarrow 90.4	-0.8	50
	LEGR	76.1 \rightarrow 75.3	-0.8	92.9 \rightarrow 92.4	-0.5	47
	CCP-AC [46]	76.2 \rightarrow 75.3	-0.9	92.9 \rightarrow 92.6	-0.3	44
	RRBP [70]	76.1 \rightarrow 73.0	-3.0	92.9 \rightarrow 91.0	-1.9	45
	C-SGD-50 [13]	75.3 \rightarrow 74.5	-0.8	92.6 \rightarrow 92.1	-0.5	45
DCP [72]	76.0 \rightarrow 74.9	-1.1	92.9 \rightarrow 92.3	-0.6	44	
MOBILENETV2	AMC [20]	71.8 \rightarrow 70.8	-1.0	\rightarrow -	-	70
	LEGR	71.8 \rightarrow 71.4	-0.4	\rightarrow -	-	70
	LEGR	71.8 \rightarrow 70.8	-1.0	\rightarrow -	-	60
	DCP [72]	70.1 \rightarrow 64.2	-5.9	\rightarrow -	-	55
	METAPRUNING [37]	72.7 \rightarrow 68.2	-4.5	\rightarrow -	-	50
	LEGR	71.8 \rightarrow 69.4	-2.4	\rightarrow -	-	50

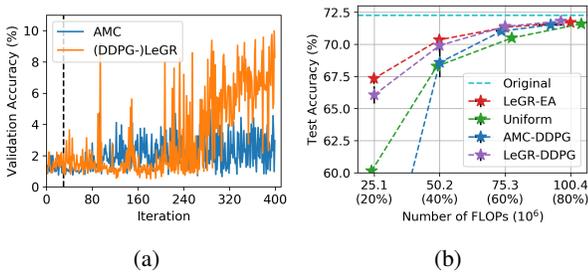


Figure 9: Comparison between searching the layer-wise filter norms and searching the layer-wise filter percentage. (a) compares the searching progress for 50% FLOP count ResNet-56 and (b) compares the final performance for ResNet-56 with various constraint levels.