

# Coreset-Based Neural Network Compression

Abhimanyu Dubey<sup>\*1</sup>, Moitreya Chatterjee<sup>\*2</sup>, and Narendra Ahuja<sup>2</sup>

<sup>1</sup> Massachusetts Institute of Technology, Cambridge MA 02139, USA  
dubeya@mit.edu

<sup>2</sup> University of Illinois at Urbana-Champaign, Champaign IL 61820, USA  
metro.smiles@gmail.com, n-ahuja@illinois.edu

**Abstract.** We propose a novel Convolutional Neural Network (CNN) compression algorithm based on coreset representations of filters. We exploit the redundancies extant in the space of CNN weights and neuronal activations (across samples) in order to obtain compression. Our method requires no retraining, is easy to implement, and obtains state-of-the-art compression performance across a wide variety of CNN architectures. Coupled with quantization and Huffman coding, we create networks that provide AlexNet-like accuracy, with a memory footprint that is  $832\times$  smaller than the original AlexNet, while also introducing significant reductions in inference time as well. Additionally these compressed networks when fine-tuned, successfully generalize to other domains as well.

## 1 Introduction

Convolutional neural networks, while immensely powerful, often are resource-intensive [35, 50, 54, 24, 26]. Popular CNN models such as AlexNet [35] and VGG-16 [50], for instance, have 61 and 138 million parameters and consume in excess of 200MB and 500MB of memory space respectively. This characteristic of deep CNN architectures reduces their portability, and poses a severe bottleneck for implementation in resource constrained environments [17]. Additionally, design choices for CNN architectures, such as network depth, filter sizes, and number of filters seem arbitrary and motivated purely by empirical performance at a particular task, permitting little room for interpretability. Moreover, the architecture design is not necessarily fully optimized for the network to be yielding a certain level of precision, making these models highly resource-inefficient.

Several prior approaches have thus sought to reduce the computational complexity of these models. Work aimed at designing efficient CNN architectures, such as Residual Networks (ResNets) [25] and DenseNets [28] have shown promise at alleviating the challenge of model complexity. These CNNs provide higher performance on classification at only a fraction of the number of parameters of their more resource intensive counterparts. However, despite being more compact, redundancies remain in such networks, leaving room for further compression.

In this work, we propose a novel method that exploits inter-filter dependencies extant in the convolutional filter banks of CNNs to compress pre-trained

---

\* Equal Contribution.

computationally intensive neural networks. Additionally we leverage neuronal activation patterns across samples to prune out irrelevant filters. Our compression pipeline consists of finely pruning the filters of every layer of the CNN based on sample activation patterns, followed by the construction of efficient filter *coreset* representations to exploit the inter-filter dependencies. Our method *does not require retraining*, is applicable to both *fully-connected and convolution layers*, and maintains classification performance similar to the uncompressed network. We display state-of-the-art compression rates on several popular CNN models, including multiple ResNets, which show increases from  $9.2\times$  to  $16.2\times$  in compression rate over prior state-of-the-art techniques. Coupled with *Deep Compression*, we are additionally able to compress other popular CNN models such as VGGNet-16 [50] and AlexNet [35] by  $238\times$  and  $55\times$  respectively. Moreover, we demonstrate the presence of filter redundancies even in highly efficient models such as SqueezeNet [29], by reducing their parameters by 50% with almost no loss in classification performance, giving us AlexNet-level precision but with  $832\times$  smaller model size, compared to the original AlexNet model. Finally, we empirically validate the generalizability of these compressed CNNs to newer domains.

In the next section, we discuss relevant prior work in this area. In Section 3, we present the details of our algorithm. This is followed by Section 4, where a discussion on the empirical evaluation of our method vis-à-vis other competing compression techniques is presented. We finally conclude in Section 5, laying out some avenues for future research in this area.

## 2 Related Work

**Network Compression:** Compressing neural networks has been a topic of active research interest lately. Prior work in this area can be grouped into three distinct categories. The first category of methods direct their attention to the construction of parameter-efficient neural network architectures. For instance, Iandola *et al.*[29] propose SqueezeNets, a neural architecture class containing the parameter efficient, fully convolutional, ‘fire’ modules. Other examples of such architectures include Residual Networks (ResNets)[25], and Densely Connected Neural Networks (DenseNets)[28], which provide higher classification performance with models much smaller than the previous state-of-the-art, using ‘skip-connections’ between layers of the network. More recent approaches have sought to adapt CNN architectures so as to make them robust to common transformations (e.g. rotation) within the data, by modifying the filter banks of a CNN [61,10] or by enforcing sparsity while training [2]. While these approaches seem to hold promise but they fail to fully exploit the inter-filter dependencies, allowing room for further compression of such networks. Meta-learning approaches attempt to decipher the optimum CNN architecture by searching over the space of a gigantic number of possible candidates. However, these techniques are prohibitively resource intensive, (needing well in excess of 400 GPUs to run), and often yield only a locally optimum architecture [46].

A second broad category of compression methods attempt to prune the unimportant network parameters. Han *et al.*[20] demonstrate an efficient pruning-retraining method, based on pruning weights by their  $\ell^p$  norms. Srinivas and Babu [53] remove individual neurons instead of weights, with impressive results. The importance of ordering filters for the purpose of pruning has also been highlighted in Yu *et al.*, He *et al.*, and Molchanov *et al.* [60,27,42]. These approaches have been modified in the works of Polyak *et al.*[45] and Luo *et al.*[41], that focus on removing weights grouped by characteristics of filters (such as norm of filter weights, etc.). Li *et al.*[40] extend the ideas of filter pruning by removing filters from a network following an ‘importance’ criterion. However, the re-training step in these algorithms is time intensive.

Finally, the third theme of compression techniques is to employ weight-approximation and information-theoretic principles for the compression of neural network parameters. An early example of such work is the approach by Denton *et al.*[9] that uses low-rank approximations to compress fully-connected layers of neural networks. However, this technique doesn’t apply to the convolution layers. Lebdev *et al.* fixes this problem and employs a low-rank decomposition approach to the full CNN to construct more efficient representations but their technique’s principal bottleneck is re-training, which we avoid [36]. Rosenfield *et al.* consider an efficient utilization of CNN filters by representing them as a linear combination of a bases set [48]. However, our algorithm, different from this line of work, is additionally also capable of introducing structure, such as sparsity, in the approximated weights resulting from the decomposition, which further aids compression. Han *et al.*[20] introduce *Deep Compression*, that uses several steps such as weight-pruning, weight-sharing, and Huffman coding to reduce neural network size. However, their algorithm requires special hardware for inference in the compressed state, making it hard to deploy the compressed networks across platforms.

Our method aims at handling the shortcomings in each of these individual themes of CNN compression. Contrary to the first category of architecture search, our method is applicable to a wide variety of models, is less resource intensive, and does not require any retraining. While we do prune filters inspired by the work of Polyak *et al.*[45] (following the second theme of compression), our criterion for filter pruning, however, is motivated by the accurate reconstruction of sample activations, instead of the magnitude of filter weights. Finally, our compression technique does not require special hardware for running inference unlike Han *et al.*[19] and scales to both fully-connected and convolutional layers, unlike the low-rank (SVD) approach by Denton *et al.*[9].

**Coresets for Point Selection:** Coresets have been widely studied in computational geometry. They were introduced first by Agarwal *et al.*[1] for approximating a set of points with a smaller set, while preserving some desired criteria, on k-means and k-median problems. Badoiu *et al.*[5] propose a coreset formulation to cluster points using a subset of the total set of points to generate the optimal solution. Har-Peled and Mazumdar [22] give an alternate solution for coresets that include points not in the original set. Feldman *et al.*[13] demonstrate that weak

coreset representations can be generated with the number of points independent of the underlying data distribution. These formulations have recently been applied to several problems within computer vision and machine learning [14,15,12], and are primarily used to approximate a set of  $n$  points in  $d$  dimensions, originating from a domain  $\mathbf{S}$ , with a smaller set of  $\tilde{n} \ll n$  points, while preserving some criterion such as similar pairwise distances. However, coresets have remained unexplored in the context of CNN compression, which constitutes a major novelty of our work.

### 3 Method

We begin with a fully-trained CNN and compress it without retraining, first by pruning out unimportant filters, followed by extraction of efficient coreset representation of these filters. Some of the major advantages of our method include: (i) Lack of retraining, therefore a major reduction in processing time, (ii) Capacity of our algorithm to significantly compress both convolutional and fully connected layers, and (iii) Ability of the compressed CNN to generalize to newer tasks.

#### 3.1 Background and Notation

An  $n$ -layered neural network can be described as a union of the parameter tensors of every layer,  $\mathcal{W} = \cup_{k=1}^n \mathbf{W}_k$ . The parameters  $\mathbf{W}_k$  of layer  $k$  have the shape  $N_k \times C_k \times h_k \times w_k$ , where  $N_k$  denotes the number of filters,  $C_k$  denotes the number of input channels of the filter (since this is typically equal to the number of filters in the previous layer,  $C_k = N_{k-1}$ ), and  $h_k$  and  $w_k$  denote the height and width of a filter. We can rewrite the parameter tensor  $\mathbf{W}_k$  as a 2D matrix  $\mathbf{W}_k$  of the shape  $N_k \times (C_k h_k w_k)$ . Next, we append the biases of the filters to  $\mathbf{W}_k$ , to make it a matrix of dimensions  $N_k \times (C_k h_k w_k + 1)$ . It is well known that using this representation of the weights and biases of a layer,  $\mathbf{W}_k$ , we can represent the output activation of any fully connected layer as a matrix product of  $\mathbf{W}_k$  with the incoming activation tensor  $\mathbf{A}_{k-1}$ . This notion can be extended to convolution layers by re-casting the matrices in an appropriate Toeplitz form [56].

The goal of compression is to obtain a compressed representation of the parameters for each layer  $\hat{\mathcal{W}} = \cup_{k=1}^n \hat{\mathbf{W}}_k$  such that it is smaller and computationally efficient, and preserves the final classification accuracy. Our approach is to construct compressed filter ‘coresets’  $\hat{\mathbf{W}}_k \in \mathbb{R}^{\hat{N}_k \times (C_k h_k w_k + 1)}$  of the parameters of each layer (where  $\hat{N}_k < N_k$ ), such that the output activations (obtained after the Toeplitz matrix multiplication), are well approximated. Ensuring that the output activations remain largely the same at every layer, post compression, ensures that the final classification performance remains largely unchanged. Since the elements of these coresets are typically linear functions of the original parameters, we will additionally require a decompression matrix  $\mathbf{D}_k \in \mathbb{R}^{\hat{N}_k \times N_k}$  to obtain an approximation to the initial set of parameters, starting with the coreset representation.

Coresets are an effective technique of approximating a large set of points with a smaller set, which need not necessarily be a part of the original set, while preserving some desirable property such as mean pairwise distances, diameter of the point set, etc. We seek to obtain a reduced matrix (coreset)  $\hat{\mathbf{W}}_k$  representation of the original filter weights  $\mathbf{W}_k$  of every layer, which we do via 3 different approaches, as described below.

### 3.2 k-Means Coresets

A first approach to constructing such a coreset would be to obtain a reduced representation of the parameter matrix that approximates the sum of distances in the space of neuronal activations of an arbitrary sample between each of the filters. Feldman *et al.*[14] demonstrate that this problem is equivalent to finding a low-rank approximation of the filter matrix. This is representable as follows:

$$\min_{\mathbf{U}'_k, \boldsymbol{\Sigma}'_k, \mathbf{V}'_k} \|\mathbf{W}_k - \mathbf{U}'_k \boldsymbol{\Sigma}'_k \mathbf{V}'_k{}^T\|_F^2 \quad (1)$$

Their formulation for constructing a compact set of  $\hat{N}_k \ll N$  points using the sum of distances criterion leads to the *k-Means Coresets*, where the coreset representation is given by the solution to the above optimization problem:

$$\hat{\mathbf{W}}_k = \mathbf{U}'_k \boldsymbol{\Sigma}'_k, \text{ with decomposition matrix: } \mathbf{D}_k = \mathbf{V}'_k{}^T \quad (2)$$

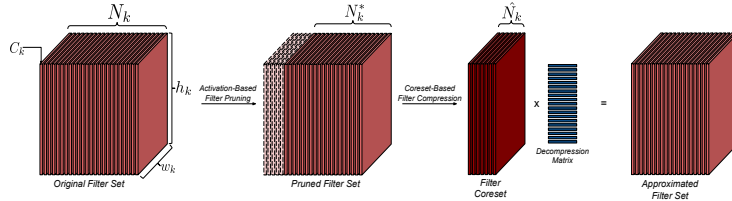
Here, the matrices  $\mathbf{U}'_k$ ,  $\boldsymbol{\Sigma}'_k$  and  $\mathbf{V}'_k{}^T$  are the  $\hat{N}_k$ -truncated versions of the matrices  $\mathbf{U}_k$ ,  $\boldsymbol{\Sigma}_k$  and  $\mathbf{V}_k{}^T$ , which satisfy the property:

$$\mathbf{W}_k = \mathbf{U}_k \boldsymbol{\Sigma}_k \mathbf{V}_k{}^T \approx \hat{\mathbf{W}}_k = \mathbf{U}'_k \boldsymbol{\Sigma}'_k \mathbf{V}'_k{}^T \quad (3)$$

$\mathbf{U}_k$ ,  $\mathbf{V}_k$  are unitary matrices, while  $\boldsymbol{\Sigma}_k$  is a diagonal matrix. Such a decomposition can be obtained using Singular Value Decomposition (SVD), where the extent of truncation is specified as an input to the algorithm. The truncation determines the amount of compression we get.

Intuitively a significant truncation, while yielding greater compression, leads to a weaker approximation of the filter weights. This also results in a weaker approximation of the output activations, manifesting itself as a drop in classification accuracy. We seek the optimum compression, across all layers, such that the classification accuracy does not deviate by more than 0.5%.

SVD for compressing neural network weights has been investigated previously in [9], however, with two key differences - (i) the naive SVD approach has been applied only to fully-connected layers of neural networks, with limited success, whereas our coreset-based formulation scales to both convolution and fully connected layers, and (ii) our method for selecting the number of components to be retained  $\hat{N}_k$  is data-dependent, based on training error obtained on random subsets of the training data, instead of an arbitrary initialization followed by retraining. However, since this decomposition does not explicitly encode any structure on the approximated weights, such as sparsity or considers the impact of activations, we build upon this formulation to create stronger coreset representations. This sets us apart from prior work, which employ simple low-rank decomposition for constructing efficient CNNs [36].



**Fig. 1.** Visual representation of compression pipeline for layer  $k$  of a neural network. Our algorithm proceeds in two steps - (i) filter pruning, and (ii) filter compression, as illustrated.

### 3.3 Structured Sparse Coresets

If we consider the previous coreset decomposition, the optimization problem can be rewritten as (subject to constraints on each of the variables  $\mathbf{U}'_k$ ,  $\mathbf{\Sigma}'_k$  and  $\mathbf{V}'_k$ ):

$$\min_{\mathbf{U}'_k, \mathbf{\Sigma}'_k, \mathbf{V}'_k} \|\mathbf{W}_k - \mathbf{U}'_k \mathbf{\Sigma}'_k \mathbf{V}'_k{}^T\|_F^2 \quad (4)$$

To induce sparsity in the obtained decomposition, Jenatton *et al.*[30] introduce a technique known as *Structured Sparse PCA*, which optimizes the following:

$$\begin{aligned} \min_{\mathbf{U}'_k, \mathbf{\Sigma}'_k, \mathbf{V}'_k} \|\mathbf{W}_k - \mathbf{U}'_k \mathbf{\Sigma}'_k \mathbf{V}'_k{}^T\|_F^2 + \lambda \cdot \|\mathbf{V}'_k\|_1, \\ \text{subject to } \|(\mathbf{U}'_k \cdot \mathbf{\Sigma}'_k)_m\|_2 = 1 \quad \forall m \in [1, \hat{N}_k] \end{aligned} \quad (5)$$

This problem can be solved by a cyclic optimization of two convex problems [30], and provides us with a decomposition that possesses structured sparsity. The motivation behind using such a formulation is to obtain a decomposition that is sparse in the number of components used, while minimizing reconstruction error. While techniques such as SPCA [32] or NMF [39] also construct representations that are sparse in the projected space, this formulation returns a decomposition that makes the approximation in the original space sparse as well, hence, *both*  $\hat{\mathbf{W}}_k$  and  $\hat{\mathbf{D}}_k$  are sparse. Moreover, this formulation allows us to discard those filters for which the corresponding column vector in  $\hat{\mathbf{D}}_k$  is a null vector, leading to further compression.

The hyper-parameters  $\hat{N}_k$  and  $\lambda$  are chosen jointly so as to obtain the maximum compression while restricting the deviation in classification performance to within 0.5% of the uncompressed network, post the compression of all layers. We observe that this technique provides much more compression than k-Means Coreset, however, this does not take into account the relative importance of the filters during reconstruction, which leads us to our final coreset formulation.

### 3.4 Activation-Weighted Coresets

Our final coreset formulation is obtained by introducing a relative importance score to every filter (based on their activation magnitudes over the training

set), while inducing sparsity. However, if we attempt to directly learn a coreset representation by minimizing the reconstruction error over all the training set activations, the resulting optimization problem will be difficult to solve, owing to the large size of the activation matrix and its degenerate nature. We thus employ an alternate formulation: for each filter  $f$  in a layer, we compute its ‘importance’  $i_k^{(f)}$  as the mean value of its activation over all training set points, normalized over all filters, in the  $k^{\text{th}}$  layer. This is given by the following:

$$i_k^{(f)} = \frac{\bar{\mathbf{A}}_k^{(f)}}{\sum_{p=1}^{N_k} \bar{\mathbf{A}}_k^{(p)}}; \quad \text{where} \quad \bar{\mathbf{A}}_k^{(f)} = \frac{1}{T} \sum_{j=1}^T \|\mathbf{A}_k^{(f)}(j)\|_F \quad (6)$$

Here,  $\mathbf{A}_k^{(f)}(j)$  is the activation of the  $f^{\text{th}}$  filter of layer  $k$ , for training sample  $j$ , and  $T$  denotes the total number of training samples. We then construct the *Importance Matrix*  $\mathbf{I}_k$  for the layer  $k$  by tiling the column vectors  $(i_k^{(f)})_{f=1}^{N_k}$ , for  $(C_k \times h_k \times w_k + 1)$  times, creating an *Importance Matrix*  $\in \mathbb{R}^{N_k \times (C_k h_k w_k + 1)}$ , where each row denotes the ‘importance’ of each filter, normalized over all filters of the current layer.

We create this form of the importance matrix with every element of a row containing identical values, since we do not want to weigh each component of a particular filter differently. Note, additionally, that we can compute this matrix in only one forward pass of the entire training set. This leads us to the following optimization problem:

$$\min_{\mathbf{U}'_k, \Sigma'_k, \mathbf{V}'_k} \|\mathbf{I}_k \odot (\mathbf{W}_k - \mathbf{U}'_k \Sigma'_k \mathbf{V}'_k{}^T)\|_F^2 \quad (7)$$

Here  $\odot$  denotes the Hadamard (elementwise) product. This problem is essentially, a weighted low-rank decomposition, studied previously by Srerbo and Jaakkola [52] and Delchambre [8] and is solved using an efficient Expectation-Maximization (EM) algorithm [52].

The intuition behind this weighted formulation is to ascribe a relative importance to the filters that contribute most to the activations in the training set (on average in the Frobenius norm sense), instead of attempting to reconstruct all activations with equal priority. Molchanov *et al.* [42] also use the notion of an importance criteria for compression but rather than using it as a weighting scheme in the optimization objective, like we do, they directly use it to prune the ‘less important’ filters. In this case as well, we compute the optimum number of components to be kept by selecting the least number of components that can be selected such that the classification accuracy is bounded within 0.5% of the original network, once the entire network has been compressed.

### 3.5 Activation-Based Filter Pruning

In related work, Liet *al.* observe that not all filters are equally important in the context of classification [40]. This motivates us to perform a pre-processing step before coreset compression, to first eliminate unimportant filters pre-emptively,

based on the mean of their activation norms over the training set. This step is essential to remove unimportant weights, since pruning out a filter in a layer can completely remove the weights corresponding to that filter, in the next layer as well, inducing greater sparsity. Using the notation from earlier we can write the size of filters  $\hat{\mathbf{W}}_k$  as:

$$\text{size}(\hat{\mathbf{W}}_k) = \hat{N}_k \times (C_k h_k w_k + 1) + N_k \hat{N}_k$$

Setting  $C_k = N_{k-1}$  (since number of outgoing activations in the previous layer is equal to the number of input channels in the next layer), and using  $N_{k-1} h_k w_k \gg N_k$ , we get:

$$\text{size}(\hat{\mathbf{W}}_k) \propto \hat{N}_k \cdot N_{k-1}$$

By layer-wise pruning of complete filters, we can hence set the number of post-pruning filters at layer  $k-1$  to be  $N_{k-1}^* < N_{k-1}$ , permitting further compression. In networks with skip-connections (e.g. ResNets),  $C_k \neq N_{k-1}$ , but it is a positive linear combination of the number of filters of the ‘‘source’’ layers of the (skip) connections, hence the proportionality still holds.

Starting from the first layer in the network, we proceed to evaluate the activation values for the entire training set, layer-by-layer. Inspired by standard ‘‘Max-Pool’’ sub-sampling techniques prevalent in modern CNNs [35,50], we approximate the response from each filter in the convolution layers (a 2D matrix) with its maximum value (a scalar). Once we have this set of pooled filter-wise activations for all samples, we compute the mean squared norm of each filter over all the training samples, and sort the filters by this value. This technique of ordering filters differentiates us from prior pruning-based techniques. We maximize the number of pruned filters, ensuring that the divergence in classification accuracy is only 0.5%, after the pruning has been carried out across all layers. Once we obtain a reduced set of filters with the crucial filters preserved, we compress this set of filters using coresets, as discussed earlier.

### 3.6 Compression Pipeline and Computational Complexity Analysis

The entire pipeline for compression can be summarized in two stages - (i) activation-based pruning, followed by (ii) coreset-based compression. The pruning procedure can be summarized in the following steps:

1. Sort the layers of the network in order of descending parameter size.
2. For each layer of the sorted network, repeat the following steps:
  - (a) Compute activations for every input in the training set, and store the maximum value for each filter activation (max-pool over spatial dimensions).
  - (b) Sort the filters in descending order of the mean value of the max-pooled activations, over the entire training set.
  - (c) Find the smallest number of filters  $N_k^*$  that can be retained while performance deviation, post the compression of all layers, is within 0.5% of original performance - using binary search.



We can see that the complexity of the individual steps are  $\mathcal{O}(n \log n)$  for the first sorting step, and  $\mathcal{O}(n \cdot (A + N_k \log N_k + A \log N_k))$  for layer-wise activation computation, filter sorting, and binary-search.  $A$  denotes the complexity to do one feed-forward operation on the entire training set. Since  $A \gg N_k \gg 1$ , the total complexity of the filter pruning is  $\mathcal{O}(n \cdot A \cdot \log \max_k N_k)$ , requiring a maximum of  $n \log \max_k N_k$  epochs of feed-forward operations, which, for most neural network architectures, we find to be much smaller than the complexity of fine-tuning.

After filter pruning, we proceed to the coreset-based compression stage. This procedure for compression can be summarized in the following steps:

For each layer in the network, starting from the shallowest, do:

1. Compute the complete decomposition according to the coreset formulation used.
2. Find the minimum number of coreset filters  $\hat{N}_k$  that can be retained while performance is within 0.5% of the network prior to coreset compression, post the compression of all layers - by searching over a random subset of the training data using binary search.

The complexity for the coreset construction set is  $\mathcal{O}(n \cdot (B + sA \log N_k^*))$ , where  $B$  is the complexity for the matrix decomposition, and  $0 \leq s \leq 1$  is the fraction of random training points used. For our experiments, we set  $s = 0.005$ . We find that for most networks,  $sA \log N_k^* > B$ , and hence the total complexity of the compression pipeline is  $\mathcal{O}(n \cdot A \cdot \log N_k \cdot (1 + s))$ . Note that post the cascading of activation-based pruning with coreset compression across all layers, the total deviation allowed in classification performance is 1% (0.5% for pruning and 0.5% for coreset compression).

## 4 Experimental Evaluation

We implement our method in PyTorch [44] and Caffe [31], and evaluate on a cluster of NVIDIA TITAN Xp and Tesla GPUs. All of our implementation and other details are available here <sup>3</sup>. For all experiments, we evaluate all 3 coreset construction techniques, as well as the impact of activation-based pruning coupled with each, and report all results together with the baseline and comparable recent work. The Activation-Based Pruning pipeline is reported as AP, while the coreset techniques are reported as - (1) k-Means Coreset (Coreset-K), (2) Structured Sparse Coreset (Coreset-S) and (3) Activation-Weighted Coreset (Coreset-A). We compare our compression performance with recent compression benchmarks, such as Fast-Food [59], SVD [9], Weight-Based Pruning [21], Deep Compression [20], memory-bounded CNNs [6], Compression Aware Training [2], etc. on a wide array of CNN architectures, including the highly efficient SqueezeNet [29].

<sup>3</sup> [https://sites.google.com/site/metrosmls/research/research-projects/compress\\_cnn](https://sites.google.com/site/metrosmls/research/research-projects/compress_cnn)

**Table 1.** Compression (Comp.) results for both AlexNet [35] and VGGNet-16 [50] trained on the ImageNet dataset, along with variation in performance with Deep Compression.

Method	AlexNet [35]				VGGNet-16 [50]		
	Acc.(%)	#Params	Comp.	#Epochs	Acc.(%)	#Params	Comp.
Baseline	57.22	61M	1×	-	68.88	138M	1×
Fastfood-32-AD [59]	58.07	30M	2×	-	-	-	-
Fastfood-16-AD [59]	57.10	17M	3.7×	-	-	-	-
Collins & Kohli [6]	55.60	15.3M	4×	-	-	-	-
Compression-Aware [2]	-	-	-	-	67.6	64.17M	2.2×
SVD [9]	55.98	12.2M	5×	540	68.85	27M	5.1×
Pruning [20]	57.23	6.8M	9×	960	68.15	15M	9.1×
Dynamic Net Surgery [18]	56.91	3.47M	17.7×	140	-	-	-
Coreset-K	56.97	9.15M	6.7×	17	68.69	15.6M	9.2×
Coreset-S	56.78	5.76M	10.5×	21	68.65	9.9M	13.9×
Coreset-A	56.82	4.97M	12.3×	23	68.01	9.2M	15.1×
AP+Coreset-K	56.51	4.02M	15.2×	26	68.56	9.81M	14×
AP+Coreset-S	56.38	<b>3.20M</b>	<b>19.1×</b>	28	67.90	<b>8.1M</b>	<b>17×</b>
AP+Coreset-A	56.48	3.68M	16.5×	27	68.16	8.7M	15.8×
<i>With Deep Compression (Comparison of Model Size)</i>							
Baseline	57.22	6.9MB	35×	-	68.70	10.77MB	49×
Coreset-K	56.80	4.17MB	49×	-	68.51	2.52MB	210×
Coreset-S	56.87	3.92MB	52×	-	68.25	2.35MB	225×
Coreset-A	57.19	4.01MB	51×	-	68.43	2.41MB	220×
AP+ Coreset-K	56.85	4.01MB	51×	-	68.02	2.28MB	232×
AP+ Coreset-S	56.70	3.85MB	53×	-	68.16	2.26MB	233×
AP+ Coreset-A	57.08	<b>3.74MB</b>	<b>55×</b>	-	68.14	<b>2.21MB</b>	<b>238×</b>

#### 4.1 LeNet-5 on MNIST

The first architecture we evaluate is the LeNet-5 network [38] on the MNIST dataset [37]. This is a popular benchmark for network compression, and high values of compression are reported by various recent work, which makes it a very competitive setup. The results for this experiment are summarized in Table 2. We can see that the coreset-based methods outperform the recent work comfortably, with a relative improvement of 18% over the existing state-of-the-art.

#### 4.2 Large-Scale ImageNet Models

The next set of experiments we perform are on the large-scale ImageNet-trained models - the very deep networks such as Residual Networks [24], AlexNet [35] and VGGNet-16 [50]. These architectures are ubiquitous for countless applied computer vision tasks [23,7], and several recent compression techniques demonstrate remarkable compression on these models which makes them an appropriate benchmark for evaluating compression performance. For these networks, we also demonstrate the impact of coupling *Deep Compression* (which involves quantization, pruning, re-training iteratively) with our method.

Table 3 summarizes the empirical evaluation on Residual Networks. We find state-of-the-art performance achieved by all three coreset methods, and a substantial increase from previous baselines as well. Even in 101-layer deep networks such as ResNet-101, we are able to obtain consistent compression, similar to the shallower ResNets. Note that this improvement is entirely on

convolutional layers, which typically have very few redundancies when compared to fully-connected layers. We additionally observe that activation-based pruning buys us significant compression, providing in essence a cascading additive effect.

Table 1 summarizes the empirical evaluation on AlexNet and VGGNet-16 networks, the two of the largest image classification networks in use today. We demonstrate substantial improvements over the state-of-the-art, by compressing AlexNet by  $19\times$ , and VGGNet-16 by  $17\times$  from their baseline sizes. When combined with Deep Compression, these ratios increase, up to  $55\times$  and  $238\times$  respectively, yielding models with a memory footprint of less than 4MB. The results additionally highlight the improvement that the activation-based pruning (AP) provides, which is most prominent in the Coreset-K and Coreset-S models.

**Table 2.** Compression (Comp.) results on LeNet-5.

Method	Top-1	Comp.
Baseline	0.97	$1\times$
Wang <i>et al.</i> [58]	0.93	$16\times$
Han <i>et al.</i> [20]	0.74	$39\times$
Guo <i>et al.</i> [18]	0.91	$108\times$
SVD [9]	0.92	$118\times$
Ullric, <i>et al.</i> [55]	0.97	$164\times$
AP+Coreset-K	0.966	$165\times$
AP+Coreset-S	0.96	$192\times$
AP+Coreset-A	0.96	<b><math>193\times</math></b>

**Table 3.** Compression results on Residual Networks. Columns Acc. and Comp. represent the Top-1 accuracy and compression factor respectively.

Method	Residual Network					
	Res-18		Res-50		Res-101	
	Acc.	Comp.	Acc.	Comp.	Acc.	Comp.
Baseline [25]	0.69	$1\times$	0.75	$1\times$	0.76	$1\times$
SVD [9]	0.69	$8\times$	0.74	$9.1\times$	0.75	$9.2\times$
Pruning [21]	0.68	$5.2\times$	0.74	$6.2\times$	0.76	$6.4\times$
N2N [4]	0.67	$9.0\times$	0.73	$8.7\times$	0.74	$8.5\times$
ThiNet [41]	-	-	0.71	$2.06\times$	-	-
ThiNet [41]	-	-	0.68	$2.95\times$	-	-
AP+Coreset-K	0.69	$13.3\times$	0.74	$14.7\times$	0.75	$15.1\times$
AP+Coreset-S	0.68	<b><math>15\times</math></b>	0.74	<b><math>15.8\times</math></b>	0.75	<b><math>16.2\times</math></b>
AP+Coreset-A	0.69	$14.2\times$	0.74	$15.6\times$	0.75	$15.8\times$

### 4.3 SqueezeNet

We evaluate our method on the highly parameter-efficient SqueezeNet architecture to evaluate if further redundancies still persist after such a compression in the architecture space and if those can be eliminated via efficient filter bank representations. We find that despite beginning with  $50\times$  less parameters than AlexNet (while providing the same performance), SqueezeNet can be compressed further (results in Table 4). Using our method, we are able to compress SqueezeNet to half its parameters, providing accuracy similar to AlexNet at  $100\times$  compression. By coupling with Deep Compression, we obtain a net compression in model size to the tune of  $16.64\times$  over the original model (or  $832\times$  from AlexNet) while maintaining classification performance.

### 4.4 Additional Observations

Further, we observe that Coreset-S and Coreset-A formulations consistently outperform Coreset-K. We surmise that large extant model redundancies tend

**Table 4.** Comparison with SqueezeNet [29] trained on the ImageNet dataset. We can compress SqueezeNet to create a model that is **832** $\times$  smaller than AlexNet [35] with the same performance.

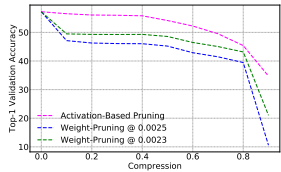
Method	Acc.(%)	Num. of Params	Ratio	Rel. to AlexNet
Baseline	57.01	1.24M	1 $\times$	50 $\times$
Coreset-K	56.83	0.73M	1.7 $\times$	85 $\times$
Coreset-S	56.92	0.65M	1.9 $\times$	95 $\times$
Coreset-A	56.94	0.61M	2 $\times$	102 $\times$
AP+ Coreset-K	56.52	0.65M	1.9 $\times$	95 $\times$
AP+ Coreset-S	56.44	0.59M	<b>2.1</b> $\times$	<b>109</b> $\times$
AP+ Coreset-A	56.80	0.60M	2 $\times$	103 $\times$
<i>With Deep Compression (Comparing Model Size)</i>				
Baseline	56.04	0.47MB	10.14 $\times$	507 $\times$
Coreset-K	56.08	0.29MB	16.1 $\times$	805 $\times$
Coreset-S	56.05	0.28MB	16.34 $\times$	817 $\times$
Coreset-A	56.03	0.29MB	16.23 $\times$	812 $\times$
AP+ Coreset-K	56.31	0.27MB	16.50 $\times$	825 $\times$
AP+ Coreset-S	56.15	0.26MB	<b>16.64</b> $\times$	<b>832</b> $\times$
AP+ Coreset-A	56.18	0.27MB	16.56 $\times$	828 $\times$

**Table 5.** LeNet-5 layer-wise compression of our method (denoted by identifiers) vis-à-vis prior work. The entries represent the fraction of parameters retained post compression.

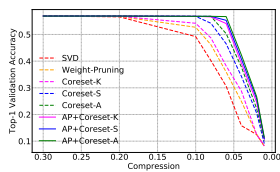
Layer	Han <i>et al.</i> [21]	Guo <i>et al.</i> [18]	K	S	A	AP+K	AP+S	AP+A
conv1	0.66	0.14	0.06	0.03	0.03	0.02	<b>0.02</b>	0.02
conv2	0.12	0.03	0.04	0.03	0.03	0.02	<b>0.02</b>	0.02
fc1	0.08	<b>0.01</b>	0.04	0.03	0.03	0.02	<b>0.01</b>	0.02
fc2	0.19	0.04	0.02	0.01	0.02	0.01	<b>0.01</b>	0.01

to benefit Coreset-A and S formulations where sparsity is explicitly enforced in the objective. Moreover, we observe that for deeper models Coreset-S tends to achieve the most compression. Table 5 shows the superior layer-wise compression achieved by our algorithm vis-à-vis state-of-the-art compression techniques on LeNet-5. The results clearly bring out the efficacy of using our compression technique, especially for convolution layers. For layer-wise compression results on other CNNs, please refer to the supplementary.

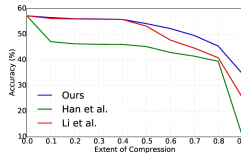
**Runtime Analysis:** We also perform a study of runtime analysis in both training and inference performance. Since we do not undertake retraining, our method is considerably faster - on our hardware, one forward pass and backward pass of AlexNet (batch size 256) takes 16ms naively, which corresponds to a total epoch training time (on ImageNet) of 2.5 minutes. We use this as a base measurement to compare the **total** training time (inclusive of the coreset operations). Table 1 describes the comparison of training times across methods. The previous state of the art method, Dynamic Net Surgery [18], requires 140 epochs (in time units) whereas our method takes *at most* 28 epochs (in time units), a significant reduction of 80%. During inference, we observe a reduction in inference time as well, which can be optimized by using efficient tensor multiplication [51]. On ResNet-50, VGGNet-16 and AlexNet, the naive (uncompressed) runtimes per epoch are: 36ms, 45ms and 8ms respectively. Our best runtimes for these



**Fig. 2.** The comparison of Activation-Based Pruning (AP) with weight-based filter pruning (without re-training), on AlexNet.



**Fig. 3.** Variation of classification performance with compression for all coreset compression techniques evaluated on AlexNet [35].



**Fig. 4.** The comparison of Activation-Based Pruning (AP) with the pruning techniques of Han *et al.* and Li *et al.* [21,40], on AlexNet.

networks (with Coreset-S) are 19ms, 21ms and 3.5ms on average, which is an average improvement of around 50%.

### 4.5 Ablation Analysis

To demonstrate the effect of individual components in our method, we perform some ablation studies as well. We first compare the effect of activation-based pruning (AP) on all coreset compression techniques on three models - AlexNet[35], VGGNet-16[50] and SqueezeNet[29], and observe that pruning benefits all methods of coreset compression, as described in Tables 1 and 4.

Next, we compare activation-based pruning with weight-based pruning, without re-training, and the pruning technique of Li *et al.* [40] for AlexNet. The results of these comparisons are summarized in Figures 2. We obtain consistently better performance at all compression ratios, substantiating the merit of data-dependent filter pruning approaches over those based on the magnitudes of filter weights.

Finally, we analyze the variation of performance with compression factor for all coreset compression techniques on the AlexNet classification model, described in Figure 3. We observe that Coreset-K (with and without AP), while stronger than SVD and Pruning approaches, worsens much more rapidly in comparison to other corresponding coreset techniques. This observation is consistent across all models. For additional results, more layer-wise compression analysis, and filter visualizations we refer the reader to the supplementary material.

### 4.6 Domain Adaptability

To measure the generalizability of our compressed models to newer tasks, we evaluate compressed models on domain adaptation benchmarks, following the experimental pipeline proposed in [41]. We evaluate the performance of the compressed CNN model VGGNet-16 [50] on target domain adaptation datasets - CUB-2011[57] and Stanford-Dogs [34], two popular datasets for fine-grained image classification. These results are summarized in Table 6. We observe that our compressed coreset models are able to provide classification performance

**Table 6.** Performance of coreset-based compression on domain-adaptation tasks.

Dataset	Model	#Params	Top-1
CUB-2011	VGG-16 Finetune(FT)	138M	72.30%
	Train from Scratch	138M	44.27%
	SVD[9] + FT	27M	53.65%
	Pruning[21] + FT	15M	57.45%
	AP+Coreset-S + FT	8.1M	<b>70.66%</b>
Stanford-Dogs	VGG-16 Finetune(FT)	138M	61.92%
	Train from Scratch	138M	27.16%
	SVD[9] + FT	27M	40.84%
	Pruning[21] + FT	15M	43.28%
	AP+Coreset-S + FT	8.1M	<b>55.91%</b>

close to the uncompressed networks, while surpassing networks compressed by other techniques. This exhibits the versatility of coreset-compressed models to domain adaptation tasks, as well.

## 5 Conclusions and Future Work

In this paper we introduce a novel technique that exploits redundancies in the space of convolutional filter weights and sample activations to reduce neural network size, using the long-existing concepts of coresets, coupled with an activation-based pooling technique. The lack of a re-training step in our algorithmic pipeline makes the implementation simple. Empirical evaluation reveals that our algorithm outperforms all other competing methods at compressing a wide array of popular CNN architectures. Our findings uncover the existence of redundancies even in the most compressed CNNs, such as SqueezeNets, which can be further exploited to improve efficiency.

Our method does not require any retraining, scales to both convolution and fully connected layers, and is extensively generalizable to different neural network models without being computationally intensive. Thus, we hope that our algorithm will serve as a valuable tool to obtain leaner and more efficient CNNs. As future work, we hope to apply our algorithm to compress other types of deep neural networks, such as Recurrent Neural Networks (RNNs) which are applicable to time-varying sequential inputs.

**Acknowledgments:** We are grateful to Prof. Ramesh Raskar for his insightful comments. MC additionally acknowledges Po-han Huang for helpful discussions and NVIDIA for providing the GPUs used for this research.

## References

1. Agarwal, P.K., Har-Peled, S., Varadarajan, K.R.: Geometric approximation via coresets. *Combinatorial and computational geometry* **52**, 1–30 (2005)

2. Alvarez, J.M., Salzmann, M.: Compression-aware training of deep networks. In: *Advances in Neural Information Processing Systems*. pp. 856–867 (2017)
3. Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Lawrence Zitnick, C., Parikh, D.: Vqa: Visual question answering. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 2425–2433 (2015)
4. Ashok, A., Rhinehart, N., Beainy, F., Kitani, K.M.: N2n learning: Network to network compression via policy gradient reinforcement learning. *arXiv preprint arXiv:1709.06030* (2017)
5. Bădoiu, M., Har-Peled, S., Indyk, P.: Approximate clustering via core-sets. In: *Proceedings of the thirty-fourth annual ACM symposium on Theory of computing*. pp. 250–257. ACM (2002)
6. Collins, M.D., Kohli, P.: Memory bounded deep convolutional networks. *arXiv preprint arXiv:1412.1442* (2014)
7. Dai, J., Qi, H., Xiong, Y., Li, Y., Zhang, G., Hu, H., Wei, Y.: Deformable convolutional networks. *CoRR*, abs/1703.06211 1(2), 3 (2017)
8. Delchambre, L.: Weighted principal component analysis: a weighted covariance eigendecomposition approach. *Monthly Notices of the Royal Astronomical Society* **446**(4), 3545–3555 (2014)
9. Denton, E.L., et al.: Exploiting linear structure within convolutional networks for efficient evaluation. In: *Advances in NIPS 2014*. pp. 1269–1277 (2014)
10. Dieleman, S., De Fauw, J., Kavukcuoglu, K.: Exploiting cyclic symmetry in convolutional neural networks. *arXiv preprint arXiv:1602.02660* (2016)
11. Dosovitskiy, A., Springenberg, J.T., Riedmiller, M., Brox, T.: Discriminative unsupervised feature learning with convolutional neural networks. In: *Advances in Neural Information Processing Systems*. pp. 766–774 (2014)
12. Dubey, A., Naik, N., Raviv, D., Sukthankar, R., Raskar, R.: Coreset-based adaptive tracking. *arXiv preprint arXiv:1511.06147* (2015)
13. Feldman, D., Monemizadeh, M., Sohler, C.: A ptas for k-means clustering based on weak coresets. In: *Proceedings of the twenty-third annual symposium on Computational geometry*. pp. 11–18. ACM (2007)
14. Feldman, D., Schmidt, M., Sohler, C.: Turning big data into tiny data: Constant-size coresets for k-means, pca and projective clustering. In: *Proceedings of the twenty-fourth annual ACM-SIAM symposium on Discrete algorithms*. pp. 1434–1453. SIAM (2013)
15. Feldman, D., Volkov, M., Rus, D.: Dimensionality reduction of massive sparse datasets using coresets. *arXiv preprint arXiv:1503.01663* (2015)
16. Girshick, R.: Fast r-cnn. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 1440–1448 (2015)
17. Gokhale, V., Jin, J., Dundar, A., Martini, B., Culurciello, E.: A 240 g-ops/s mobile coprocessor for deep neural networks. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*. pp. 682–687 (2014)
18. Guo, Y., Yao, A., Chen, Y.: Dynamic network surgery for efficient dnns. In: *Advances In Neural Information Processing Systems*. pp. 1379–1387 (2016)
19. Han, S., Liu, X., Mao, H., Pu, J., Pedram, A., Horowitz, M.A., Dally, W.J.: Eie: efficient inference engine on compressed deep neural network. In: *Proceedings of the 43rd International Symposium on Computer Architecture*. pp. 243–254. IEEE Press (2016)
20. Han, S., Mao, H., Dally, W.J.: Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding (2015)

21. Han, S., Pool, J., Tran, J., Dally, W.: Learning both weights and connections for efficient neural network. In: *Advances in Neural Information Processing Systems*. pp. 1135–1143 (2015)
22. Har-Peled, S., Mazumdar, S.: On coresets for k-means and k-median clustering. *ACM Symposium on Theory of Computing* (2004)
23. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: *Computer Vision (ICCV), 2017 IEEE International Conference on*. pp. 2980–2988. IEEE (2017)
24. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385* (2015)
25. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 770–778 (2016)
26. He, K., Zhang, X., Ren, S., Sun, J.: Identity mappings in deep residual networks. *arXiv preprint arXiv:1603.05027* (2016)
27. He, Y., et al.: Channel pruning for accelerating very deep nns. In: *ICCV 2017*
28. Huang, G., Liu, Z., Weinberger, K.Q., van der Maaten, L.: Densely connected convolutional networks. *arXiv preprint arXiv:1608.06993* (2016)
29. Iandola, F.N., Moskewicz, M.W., Ashraf, K., Han, S., Dally, W.J., Keutzer, K.: Squeezenet: Alexnet-level accuracy with 50x fewer parameters and < 1mb model size. *arXiv preprint arXiv:1602.07360* (2016)
30. Jenatton, R., Obozinski, G., Bach, F.: Structured sparse principal component analysis. In: *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*. pp. 366–373 (2010)
31. Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., Darrell, T.: Caffe: Convolutional architecture for fast feature embedding. In: *Proceedings of the 22nd ACM international conference on Multimedia*. pp. 675–678. ACM (2014)
32. Jolliffe, I.T., Trendafilov, N.T., Uddin, M.: A modified principal component technique based on the lasso. *Journal of computational and Graphical Statistics* **12**(3), 531–547 (2003)
33. Karpathy, A., Fei-Fei, L.: Deep visual-semantic alignments for generating image descriptions. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 3128–3137 (2015)
34. Khosla, A., Jayadevaprakash, N., Yao, B., Li, F.F.: Novel dataset for fine-grained image categorization: Stanford dogs
35. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: *Advances in neural information processing systems*. pp. 1097–1105 (2012)
36. Lebedev, V., Ganin, Y., Rakhuba, M., Oseledets, I., Lempitsky, V.: Speeding-up convolutional neural networks using fine-tuned cp-decomposition. *arXiv preprint arXiv:1412.6553* (2014)
37. LeCun, Y.: The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/> (1998)
38. LeCun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. *Proceedings of the IEEE* **86**(11), 2278–2324 (1998)
39. Lee, D.D., Seung, H.S.: Learning the parts of objects by non-negative matrix factorization. *Nature* **401**(6755), 788–791 (1999)
40. Li, H., et al.: Pruning filters for efficient convnets. *ICLR 2017*
41. Luo, J.H., Wu, J., Lin, W.: Thinet: A filter level pruning method for deep neural network compression. *arXiv preprint arXiv:1707.06342* (2017)



42. Molchanov, P., et al.: Pruning convolutional neural networks for resource efficient transfer learning. ICLR 2017
43. Nam, H., Han, B.: Learning multi-domain convolutional neural networks for visual tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4293–4302 (2016)
44. Paszke, A., Gross, S., Chintala, S.: Pytorch (2017)
45. Polyak, A., Wolf, L.: Channel-level acceleration of deep face representations. IEEE Access **3**, 2163–2175 (2015)
46. Real, E., Moore, S., Selle, A., Saxena, S., Suematsu, Y.L., Le, Q., Kurakin, A.: Large-scale evolution of image classifiers. arXiv preprint arXiv:1703.01041 (2017)
47. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: Advances in neural information processing systems. pp. 91–99 (2015)
48. Rosenfeld, A., Tsotsos, J.K.: Incremental learning through deep adaptation. arXiv preprint arXiv:1705.04228 (2017)
49. Shin, H.C., Roth, H.R., Gao, M., Lu, L., Xu, Z., Nogues, I., Yao, J., Mollura, D., Summers, R.M.: Deep convolutional neural networks for computer-aided detection: Cnn architectures, dataset characteristics and transfer learning. IEEE transactions on medical imaging **35**(5), 1285–1298 (2016)
50. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 (2014)
51. Solomonik, E.: Provably efficient algorithms for numerical tensor algebra. University of California, Berkeley (2014)
52. Srebro, N., Jaakkola, T.: Weighted low-rank approximations. In: Proceedings of the 20th International Conference on Machine Learning (ICML-03). pp. 720–727 (2003)
53. Srinivas, S., Babu, R.V.: Data-free parameter pruning for deep neural networks. arXiv preprint arXiv:1507.06149 (2015)
54. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1–9 (2015)
55. Ullrich, K., Meeds, E., Welling, M.: Soft weight-sharing for neural network compression. arXiv preprint arXiv:1702.04008 (2017)
56. Vasudevan, A., Anderson, A., Gregg, D.: Parallel multi channel convolution using general matrix multiplication. arXiv preprint arXiv:1704.04428 (2017)
57. Wah, C., Branson, S., Welinder, P., Perona, P., Belongie, S.: The caltech-ucsd birds-200-2011 dataset (2011)
58. Wang, S., Cai, H., Bilmes, J., Noble, W.: Training compressed fully-connected networks with a density-diversity penalty (2016)
59. Yang, Z., Moczulski, M., Denil, M., de Freitas, N., Smola, A., Song, L., Wang, Z.: Deep fried convnets. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1476–1483 (2015)
60. Yu, R., et al.: Nisp: Pruning networks using neuron importance score propagation. arXiv preprint arXiv:1711.05908 (2017)
61. Zhai, S., Cheng, Y., Zhang, Z.M., Lu, W.: Doubly convolutional neural networks. In: Advances in neural information processing systems. pp. 1082–1090 (2016)