

# TRAINING CNNs WITH LOW-RANK FILTERS FOR EFFICIENT IMAGE CLASSIFICATION

Yani Ioannou<sup>1</sup>, Duncan Robertson<sup>2</sup>, Jamie Shotton<sup>2</sup>, Roberto Cipolla<sup>1</sup> & Antonio Criminisi<sup>2</sup>

<sup>1</sup>University of Cambridge, <sup>2</sup>Microsoft Research

{yai20,rc10001}@cam.ac.uk, {a-durobe,jamiesho,antcrim}@microsoft.com

## ABSTRACT

We propose a new method for creating computationally efficient convolutional neural networks (CNNs) by using low-rank representations of convolutional filters. Rather than approximating filters in previously-trained networks with more efficient versions, we learn a set of small basis filters from scratch; during training, the network learns to combine these basis filters into more complex filters that are discriminative for image classification. To train such networks, a novel weight initialization scheme is used. This allows effective initialization of connection weights in convolutional layers composed of groups of differently-shaped filters. We validate our approach by applying it to several existing CNN architectures and training these networks from scratch using the CIFAR, ILSVRC and MIT Places datasets. Our results show similar or higher accuracy than conventional CNNs with much less compute. Applying our method to an improved version of VGG-11 network using global max-pooling, we achieve comparable validation accuracy using 41% less compute and only 24% of the original VGG-11 model parameters; another variant of our method gives a 1 percentage point *increase* in accuracy over our improved VGG-11 model, giving a top-5 *center-crop* validation accuracy of 89.7% while reducing computation by 16% relative to the original VGG-11 model. Applying our method to the GoogLeNet architecture for ILSVRC, we achieved comparable accuracy with 26% less compute and 41% fewer model parameters. Applying our method to a near state-of-the-art network for CIFAR, we achieved comparable accuracy with 46% less compute and 55% fewer parameters.

## 1 INTRODUCTION

Convolutional neural networks (CNNs) have been used increasingly successfully to solve challenging computer vision problems such as image classification (Krizhevsky et al., 2012), object detection (Ren et al., 2015), and human pose estimation (Tompson et al., 2015). However, recent improvements in recognition accuracy have come at the expense of increased model size and computational complexity. These costs can be prohibitive for deployment on low-power devices, or for fast analysis of videos and volumetric medical images.

One promising aspect of CNNs, from a memory and computational efficiency standpoint, is their use of *convolutional filters* (a.k.a. kernels). Such filters usually have limited spatial extent and their learned weights are shared across the image spatial domain to provide translation invariance (Fukushima, 1980; LeCun et al., 1998). Thus, as illustrated in Fig. 1, in comparison with fully connected network layers (Fig. 1a), convolutional layers have a much sparser connection structure and use fewer parameters (Fig. 1b). This leads to faster training and test, better generalization, and higher accuracy.

This paper focuses on reducing the computational complexity of the convolutional layers of CNNs by further sparsifying their connection structures. Specifically, we show that by representing convolutional filters using a basis space comprising groups of filters of different spatial dimensions (examples shown in Fig. 1c and d), we can significantly reduce the computational complexity of existing state-of-the-art CNNs without compromising classification accuracy.

Our contributions include a novel method of learning a set of small basis filters that are combined to represent larger filters efficiently. Rather than approximating previously trained networks, we

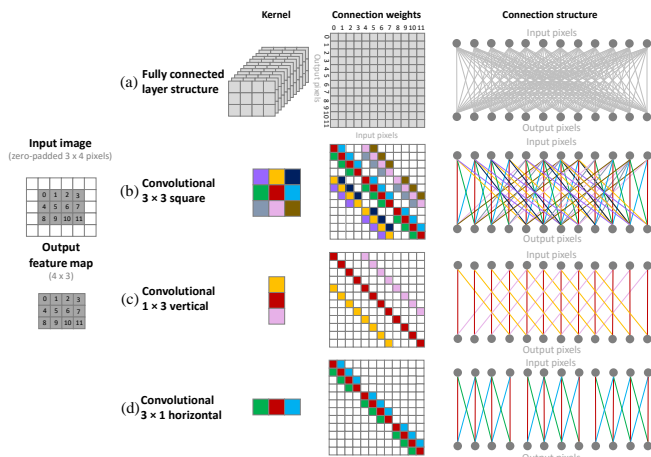


Figure 1: **Network connection structure for convolutional layers.** An input image is transformed in one layer of a neural network into an output image of same size. Connection weight maps show pairwise dependencies between input and output pixels. In (a), each node is connected to all input pixels. For (b,c,d), output pixels depend only on a subset of input pixels (shared weights are represented by unique colours). Note that sparsity increases from (a) to (d), opening up potentially more efficient implementation.

train networks *from scratch* and show that our convolutional layer representation can improve both efficiency and classification accuracy. We further describe how to initialize connection weights effectively for training networks with composite convolutional layers containing groups of differently-shaped filters, which we found to be of critical importance to our training method.

### 1.1 RELATED WORK

There has been much previous work on increasing the test-time efficiency of CNNs. Some promising approaches work by making use of more hardware-efficient representations. For example [Gupta et al. \(2015\)](#) and [Vanhoucke et al. \(2011\)](#) achieve training- and test-time compute savings by further quantization of network weights that were originally represented as 32 bit floating point numbers. However, more relevant to our work are approaches that depend on new network connection structures, efficient approximations of previously trained networks, and learning low rank filters.

**Efficient Network Connection Structures.** There has been shown to be significant redundancy in the trained weights of CNNs ([Denil et al., 2013](#)). [LeCun et al. \(1989\)](#) suggest a method of pruning unimportant connections within networks. However this requires repeated network re-training and may be infeasible for modern, state-of-the-art CNNs requiring weeks of training time. [Lin et al. \(2013\)](#) show that the geometric increase in the number and dimensions of filters with deeper networks can be managed using low-dimensional embeddings. The same authors show that global average-pooling may be used to decrease model size in networks with fully connected layers. [Simonyan & Zisserman \(2014\)](#) show that stacked filters with small spatial dimensions (*e.g.*  $3 \times 3$ ), can operate on the effective receptive field of larger filters (*e.g.*  $5 \times 5$ ) with less computational complexity.

**Low-Rank Filter Approximations.** [Rigamonti et al. \(2013\)](#) approximate *previously trained* CNNs with low-rank filters for the semantic segmentation of curvilinear structures within volumetric medical imagery. They discuss two approaches: enforcing an  $L_1$ -based regularization to learn approximately low rank filters, which are later truncated to enforce a strict rank, and approximating a set of pre-learned filters with a tensor-decomposition into many rank-1 filters. Neither approach learns low rank filters directly, and indeed the second approach proved the more successful.

The work of [Jaderberg et al. \(2014\)](#) also approximates the existing filters of previously trained networks. They find separable 1D filters through an optimization minimizing the reconstruction error of the already learned full rank filters. They achieve a  $4.5\times$  speed-up with a loss of accuracy of 1% in a text recognition problem. However since the method is demonstrated only on text recognition, it is not clear how well it would scale to larger data sets or more challenging problems. A key insight

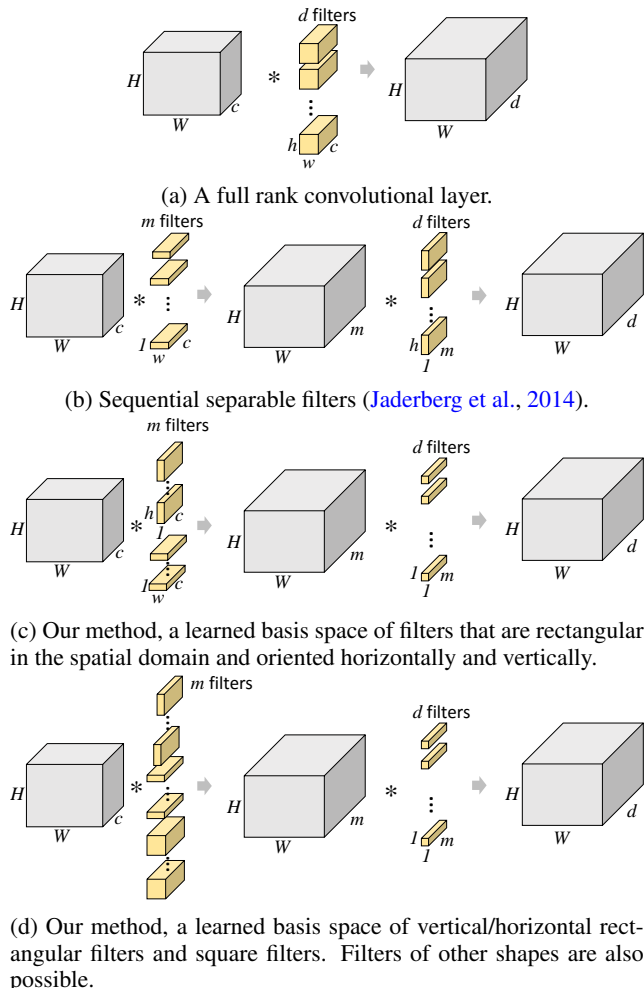


Figure 2: **Methods of using low-rank filters in CNNs.** The activation function is not shown, coming after the last layer in each configuration.

of the paper is that filters can be represented by low rank approximations not only in the spatial domain but also in the channel domain.

Both of these methods show that, at least for their respective applications, low rank approximations of full-rank filters learned in convolutional networks can increase test time efficiency significantly. However, being approximations of pre-trained networks, they are unlikely to improve test accuracy, and can only increase the computational requirements during training.

**Learning Separable Filters.** Mamalet & Garcia (2012) propose training networks with separable filters on the task of digit recognition with the MNIST dataset. They train networks with *sequential* convolutional layers of horizontal and vertical 1D filters, achieving a speed-up factor of  $1.6\times$ , but with a relative increase in test error of 13% (1.45% v.s. 1.28%). Our approach generalizes this, allowing both horizontal and vertical 1D filters (and other shapes too) at the same layer and avoiding issues with ordering. We also demonstrate a decrease in error and on more challenging datasets.

## 2 USING LOW-RANK FILTERS IN CNNs

### 2.1 CONVOLUTIONAL FILTERS

The convolutional layers of a CNN produce output ‘images’ (usually called *feature maps*) by convolving input images with one or more learned filters. In a typical convolutional layer, as illustrated

in Fig. 2a, a  $c$ -channel input image of size  $H \times W$  pixels is convolved with  $d$  filters of size  $h \times w \times c$  to create a  $d$ -channel output image. Each filter is represented by  $hwc$  independent weights. Therefore the computational complexity for the convolution of the filter with a  $c$ -channel input image is  $\mathcal{O}(dwhc)$  (per pixel in the output feature map).

In what follows, we describe schemes for modifying the architecture of the convolutional layers so as to reduce computational complexity. The idea is to replace full-rank convolutional layers with modified versions that represent the same number of filters by linear combinations of basis vectors, *i.e.* as lower rank representations of the full rank originals.

## 2.2 SEQUENTIAL SEPARABLE FILTERS

An existing scheme for reducing the computational complexity of convolutional layers (Jaderberg et al., 2014) is to replace each one with a sequence of two regular convolutional layers but with filters that are rectangular in the spatial domain, as shown in Fig 2b. The first convolutional layer has  $m$  filters of size  $w \times 1 \times c$ , producing an output feature map with  $m$  channels. The second convolutional layer has  $d$  filters of size  $1 \times h \times m$ , producing an output feature map with  $d$  channels. By this means the full rank original convolutional filter bank is represented by a low rank approximation formed from a linear combination of a set of separable  $w \times h$  basis filters. The computational complexity of this scheme is  $\mathcal{O}(mcw)$  for the first layer of horizontal filters and  $\mathcal{O}(dmh)$  for the second layer of vertical filters, with a total of  $\mathcal{O}(m(cw + dh))$ .

Note that Jaderberg et al. (2014) use this scheme to approximate existing full rank filters belonging to previously trained networks using a retrospective fitting step. In this work, by contrast, we *train* networks containing convolutional layers with this architecture from scratch. In effect, we learn the separable basis filters and their combination weights simultaneously during network training.

## 2.3 FILTERS AS LINEAR COMBINATIONS OF BASES

In this work we introduce another scheme for reducing convolutional layer complexity. This works by representing convolutional filters as linear combinations of basis filters as illustrated in Fig. 2c. This scheme uses *composite layers* comprising several sets of filters where the filters in each set have different spatial dimensions (see Fig. 5). The outputs of these basis filters may be combined in a subsequent layer containing filters with spatial dimensions  $1 \times 1$ .

This is illustrated in Fig. 2c. Here, our composite layer contains horizontal  $w \times 1$  and vertical  $1 \times h$  filters, the outputs of which are concatenated in the channel dimension, resulting in an intermediate  $m$ -channel feature map. These filter responses are then linearly combined by the next layer of  $d$   $1 \times 1$  filters to give a  $d$ -channel output feature map. In this case, the filters are applied on the input feature map with  $c$  channels and followed by a set of  $m$   $1 \times 1$  filters over the  $m$  output channels of the basis filters. If the number of horizontal and vertical filters is the same, the computational complexity is  $\mathcal{O}(m(wc/2 + hc/2 + d))$ . Interestingly, the configuration of Fig. 2c gives rise to linear combinations of horizontal and vertical filters that are cross-shaped in the spatial domain. This is illustrated in Fig. 3 for filters learned in the first convolutional layer of the ‘vgg-gmp-lr-join’ model that is described in the Results section when it is trained using ILSVRC dataset.

Note that, in general, more than two different sizes of basis filter might be used in the composite layer. For example, Fig. 2d shows a combination of three sets of filters with spatial dimensions  $w \times 1$ ,  $1 \times h$ , and  $w \times h$ . Also note that an interesting option is to omit the  $1 \times 1$  linear combination layer and instead allow the connection weights in a subsequent network layer to learn to combine the basis filters of the preceding layer (despite any intermediate non-linearity, *e.g.* ReLUs). This possibility is explored in practice in the Results section.

In that our method uses a combination of filters in a composite layer, it is similar to the ‘GoogLeNet’ of Szegedy et al. (2014) which uses ‘inception’ modules comprising several (square) filters of different sizes ranging from  $1 \times 1$  to  $5 \times 5$ . In our case, however, we are implicitly learning linear combinations of less computationally expensive filters with different orientations (*e.g.*  $3 \times 1$  and  $1 \times 3$  filters), rather than combinations of filters of different sizes. Amongst networks with similar computational requirements, GoogLeNet is one of the most accurate for large scale image classification tasks (see Fig. 4), partly due to the use of heterogeneous filters in the inception modules, but also the use of low-dimensional embeddings and global pooling.

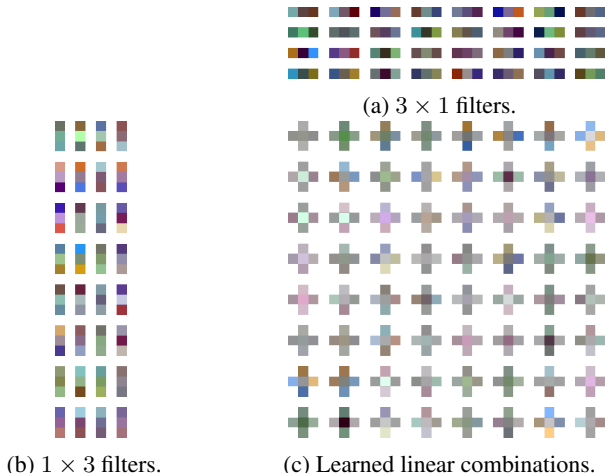


Figure 3: **Learned Cross-Shaped Filters.** The cross-shaped filters (c) learned as weighted linear combination of (b)  $1 \times 3$  and (a)  $3 \times 1$  basis filters in the first convolutional layer of the the ‘vgg-gmp-lr-join’ model trained using the ILSVRC dataset.

### 3 TRAINING CNNs WITH MIXED-SHAPE LOW-RANK FILTERS

To determine the standard deviations to be used for weight initialization, we use an approach similar to that described by [Glorot & Bengio \(2010\)](#) (with the adaptation described by [He et al. \(2015\)](#) for layers followed by a ReLU). In Appendix A, we show the details of our derivation, generalizing the approach of [He et al. \(2015\)](#) to the initialization of ‘composite’ layers comprising several groups of filters of different spatial dimensions (see Appendix A, Fig. 5). This is one of the main contributions of this work.

We find that a composite layer of heterogeneously-shaped filter groups, where each filter group  $i$  has  $w^{[i]}h^{[i]}d^{[i]}$  outgoing connections should be initialized as if it is a single layer with  $\hat{n} = \sum w^{[i]}h^{[i]}d^{[i]}$ . Thus in the case of a ReLU non-linearity, we find that such a composite layer should be initialized with a zero-mean Gaussian distribution with standard deviation:

$$\sigma = \sqrt{\frac{2}{\sum w^{[i]}h^{[i]}d^{[i]}}}. \quad (1)$$

### 4 RESULTS AND COMPARISONS

To validate our approach, we show that we can replace the filters used in existing state-of-the-art network architectures with low-rank representations as described above to reduce computational complexity without reducing accuracy. Here we characterize the computational complexity of a CNN using the number of multiply accumulate operations required for a forward pass (which depends on the size of the filters in each convolutional layer as well as the input image size and stride). However, we have observed strong correlation between multiply-accumulate counts and runtime for both CPU and GPU implementations of the networks described here (as shown in Appendix B, Fig. 6). Note that the Caffe timings differ more for the initial convolutional layers where the input sizes are much smaller (3-channels) as BLAS is less efficient for the relatively small matrices being multiplied.

**Methodology.** We augment our training set with randomly cropped and mirrored images, but do not use any scale or photometric augmentation, or over-sampling. This allows us to compare the efficiency of different network architectures without having to factor in the computational cost of the various augmentation methods used elsewhere. During training, for every model except GoogLeNet, we adjust the learning rate according to the schedule  $\gamma_t = \gamma_0(1 + \gamma_0\lambda t)^{-1}$ , where  $\gamma_0$ ,  $\gamma_t$  and  $\lambda$  are the initial learning rate, learning rate at iteration  $t$ , and weight decay respectively ([Bottou, 2012](#)).

When the validation accuracy levels off we manually reduce the learning rate by further factors of 10 until the validation accuracy no longer increases. Unless otherwise indicated, aside from changing the standard deviation of the normally distributed weight initialization, as explained in §3, we used the standard hyper-parameters for each given model. Our results use no test-time augmentation.

#### 4.1 VGG-11 ARCHITECTURES FOR ILSVRC OBJECT CLASSIFICATION AND MIT PLACES SCENE CLASSIFICATION

We evaluated classification accuracy of the VGG-11 based architectures using two datasets, ImageNet Large Scale Visual Recognition Challenge 2012 (‘ILSVRC’) and MIT Places. The ILSVRC dataset comprises 1.2M training images of 1000 object classes, commonly evaluated by top-1 and top-5 accuracy on the 50K image validation set. The MIT Places dataset comprises 2.4M training images from 205 scene classes, evaluated with top-1 and top-5 accuracy on the 20K image validation set.

VGG-11 (‘VGG-A’) is an 11-layer convolutional network introduced by [Simonyan & Zisserman \(2014\)](#). It is in the same family of network architectures used by [Simonyan & Zisserman \(2014\)](#); [He et al. \(2015\)](#) to obtain the state-of-the-art accuracy for ILSVRC, but uses fewer convolutional layers and therefore fits on a single GPU during training. During training of our VGG-11 based models, we used the standard hyperparameters as detailed by [Simonyan & Zisserman \(2014\)](#) and the initialization of [He et al. \(2015\)](#).

In what follows, we compare the accuracy of a number of different network architectures detailed in Appendix E, Table 6. Results for ILSVRC are given in Table 1, and plotted in Fig. 4. Results for MIT Places are given in Table 2, and plotted in Fig. 9.

**Baseline (Global Max Pooling).** Compared to the version of the network described in ([Simonyan & Zisserman, 2014](#)), we use a variant that replaces the final  $2 \times 2$  max pooling layer before the first fully connected layer with a global max pooling operation, similar to the global average pooling used by [Lin et al. \(2013\)](#); [Szegedy et al. \(2014\)](#). We evaluated the accuracy of the baseline VGG-11 network with global max-pooling (**vgg-gmp**) and without (**vgg-11**) on the two datasets. We trained these networks at stride 1 on the ILSVRC dataset and at stride 2 on the larger MIT Places dataset. This globally max-pooled variant of VGG-11 uses over 75% fewer parameters than the original network and gives consistently better accuracy – almost 3 percentage points lower top-5 error on ILSVRC than the baseline VGG-11 network on ILSVRC (see Table 1). We used this network as the baseline for the rest of our experiments.

**Separable Filters.** To evaluate the separable filter approach described in §2.2 (illustrated in Fig. 2b), we replaced each convolutional layer in VGG-11 with a sequence of two layers, the first containing horizontally oriented  $1 \times 3$  filters and the second containing vertically oriented  $3 \times 1$  filters (**vgg-gmp-sf**). These filters applied in sequence represent  $3 \times 3$  kernels using a low dimensional basis space. Unlike [Jaderberg et al. \(2014\)](#), we trained this network from scratch instead of approximating the full-rank filters in a previously trained network. Compared to the original VGG-11 network, the separable filter version requires approximately 14% less compute. Results are shown in Table 1 for ILSVRC and Table 2 for MIT Places. Accuracy for this network is approx. 0.8% lower than that of the baseline vgg-11-gmp network for ILSVRC and broadly comparable for MIT Places. This approach does not give such a significant reduction in computational complexity as what follows, but it is nonetheless interesting that separable filters are capable of achieving quite high classification accuracy on such challenging tasks.

**Simple Horizontal/Vertical Basis.** To demonstrate the efficacy of the simple low rank filter representation illustrated in Fig. 2c, we created a new network architecture (**vgg-gmp-lr-join**) by replacing each of the convolutional layers in VGG-11 (original filter dimensions were  $3 \times 3$ ) with a sequence of two layers. The first layer comprises half  $1 \times 3$  filters and half  $3 \times 1$  filters whilst the second layer comprises the same number of  $1 \times 1$  filters. The resulting network is approximately 49% faster than the original and yet it gives broadly comparable accuracy (within 1 percentage point) for both the ILSVRC and MIT Places datasets.

**Full-Rank Mixture.** An interesting question concerns the impact on accuracy of combining a small proportion of 33 filters with the 13 and 31 lters used in vgg-gmp-lr-join. To answer this ques-

tion, we trained a network, **vgg-gmp-lr-join-wfull**, with a mixture of 25%  $3 \times 3$  and 75%  $1 \times 3$  and  $3 \times 1$  filters, while preserving the total number of filters of the baseline network (as illustrated in Fig. 2d). This network was significantly more accurate than both ‘vgg-gmp-lr-join’ and the baseline, with a top-5 center crop accuracy of 89.7% on ILSVRC, with a computational savings of approx. 16% over our baseline. We note that the accuracy is approx. 1 percentage point higher than GoogLeNet.

**Implicitly Learned Combinations.** In addition, we try a network similar to vgg-gmp-lr-join but without the  $1 \times 1$  convolutional layer (as shown in Fig. 2c) used to sum the contributions of  $3 \times 1$  and  $1 \times 3$  filters (**vgg-gmp-lr**). Interestingly, because of the elimination of the extra  $1 \times 1$  layers, this gives an additional compute saving such that this model is only  $1/3^{\text{rd}}$  of the compute of our baseline, with no reduction in accuracy. This seems to be a consequence of the fact that the subsequent convolutional layer is itself capable of learning effective combinations of filter responses even after the intermediate ReLU non-linearity.

We also trained such a network with double the number of convolutional filters (**vgg-gmp-lr-2x**), *i.e.* with an equal number of  $1 \times 3$  and  $3 \times 1$  filters, or  $2c$  filters as shown in Fig. 2c. We found this to increase accuracy further (88.9% Top-5 on ILSVRC) while still being approximately 58% faster than our baseline network.

**Low-Dimensional Embeddings.** We attempted to reduce the computational complexity of our ‘gmp-lr’ network further in the **vgg-gmp-lr-lde** network by using a stride of 2 in the first convolutional layer, and adding low-dimensional embeddings, as in Lin et al. (2013); Szegedy et al. (2014). We reduced the number of output channels by half after each convolutional layer using  $1 \times 1$  convolutional layers, as detailed in Appendix E, Table 6. While this reduces computation significantly, by approx. 86% compared to our baseline, we saw a decrease in top-5 accuracy on ILSVRC of 1.2 percentage points. We do note however, that this network remains 2.5 percentage points more accurate than the original VGG-11 network, but is 87% faster.

Network	Stride	Multiple-Acc. $\times 10^9$	Param. $\times 10^7$	Top-1 Acc.	Top-5 Acc.
vgg-11	1	7.61	13.29	0.649	0.862
gmp	1	7.51	3.22	0.685	0.887
gmp-sf	1	6.53	2.97	0.673	0.879
gmp-lr-join-wfull	1	6.34	3.72	<b>0.704</b>	<b>0.897</b>
gmp-lr-join	1	3.85	2.73	0.675	0.880
gmp-lr-2x	1	3.14	3.13	0.693	0.889
gmp-lr	1	2.52	<b>2.61</b>	0.676	0.880
gmp-lr-lde	2	<b>1.02</b>	2.64	0.667	0.875

Table 1: **VGG ILSVRC Results.** Accuracy, multiply-accumulate count, and number of parameters for the baseline VGG-11 network (both with and without global max pooling) and more efficient versions created by the methods described in this paper.

Network	Stride	Multiple-Acc. $\times 10^8$	Param. $\times 10^7$	Top-1 Acc.	Top-5 Acc.
gmp	2	18.77	3.22	<b>0.526</b>	<b>0.830</b>
gmp-sf	2	16.57	13.03	0.517	0.824
gmp-lr-join	2	9.64	2.73	0.512	0.821
gmp-lr	2	<b>6.30</b>	<b>2.61</b>	0.520	0.825

Table 2: **MIT Places Results.** Accuracy, multiply-accumulate operations, and number of parameters for the baseline ‘vgg-11-gmp’ network, separable filter network as described by Jaderberg et al. (2014), and more efficient models created by the methods described in this paper. All networks were trained at stride 2 for the MIT Places dataset.

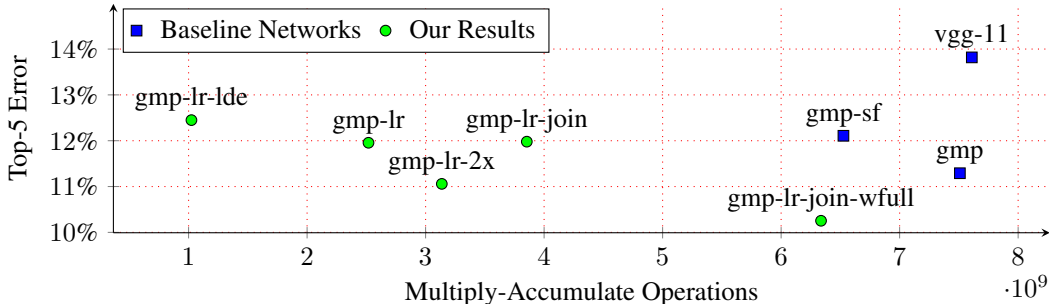


Figure 4: **VGG ILSVRC Results.** Multiply-accumulate operations v.s. top-5 error for VGG-derived models on ILSVRC object classification dataset, the most efficient networks are closer to the origin. Our models are significantly faster than the baseline network, in the case of ‘gmp-lr-2x’ by a factor of almost 60%, while slightly lowering error. Note that the ‘gmp-lr’ and ‘gmp-lr-join’ networks have the same accuracy, showing that an explicit linear combination layer may be unnecessary.

#### 4.2 GOOGLeNET FOR ILSVRC OBJECT CLASSIFICATION

GoogLeNet, introduced by Szegedy et al. (2014), is the most efficient network for ILSVRC, getting close to state-of-the-art results with a fraction of the compute and model size of even VGG-11. The GoogLeNet inception module is a composite layer of 5 homogeneously-shaped filters,  $1 \times 1$ ,  $3 \times 3$ ,  $5 \times 5$ , and the output of a  $3 \times 3$  average pooling operations. All of these are concatenated and used as input for successive layers.

For the **googlenet-lr** network, within only the inception modules we replaced each the  $3 \times 3$  filters with low-rank  $3 \times 1$  and  $1 \times 3$  filters, and replaced the layer of  $5 \times 5$  filters with a set of low-rank  $5 \times 1$  and  $1 \times 5$  filters. For the **googlenet-lr-conv1** network, we similarly replaced the first and second layer convolutional layers with  $7 \times 1 / 1 \times 7$  and  $3 \times 1 / 1 \times 3$  layers respectively.

Results are shown in Table 3. Due to the intermediate losses used for training, which contain the only fully-connected layers in GoogLeNet, test time model size is significantly smaller than training time model size. Table 3 also reports test time model size. The low-rank network delivers comparable classification accuracy using 26% less compute. No other networks produce comparable accuracy within an order of magnitude of compute. We note that although the Caffe pre-trained GoogLeNet model (Jia et al., 2014) has a top-5 accuracy of 0.889, our training of the same network using the given model definition, including the hyper-parameters and training schedule, but a different random initialization had a top-5 accuracy of 0.883.

Network	Multiple-Acc. $\times 10^9$	Test Param. $\times 10^6$	Top-1 Acc.	Top-5 Acc.
GoogLeNet	1.59	5.97	<b>0.677</b>	<b>0.883</b>
lr	1.18	3.50	0.673	0.880
lr-conv1	<b>0.84</b>	<b>3.42</b>	0.659	0.870

Table 3: **GoogLeNet ILSVRC Results.** Accuracy, multiply-accumulate count, and number of parameters for the baseline GoogLeNet network and more efficient versions created by the methods described in this paper.

#### 4.3 NETWORK-IN-NETWORK FOR CIFAR-10 OBJECT CLASSIFICATION

The CIFAR-10 dataset consists of 60,000  $32 \times 32$  images in 10 classes, with 6000 images per class. This is split into standard sets of 50,000 training images, and 10,000 test images (Krizhevsky, 2009). As a baseline for the CIFAR-10 dataset, we used the Network in Network architecture (Lin et al., 2013), which has a published test-set error of 8.81%. We also used random crops during training, with which the network has an error of 8.1%. Like most state of the art CIFAR results, this was with ZCA pre-processed training and test data (Goodfellow et al., 2013), training time mirror augmentation and random sub-crops. The results of our CIFAR experiments are listed in Table 4 and plotted in Fig. 11.



Network	Multiple-Acc. $\times 10^8$	Param. $\times 10^5$	Accuracy
NiN	1.93	9.67	<b>0.9188</b>
nin-c3	1.43	7.74	0.9186
nin-c3-lr	<b>1.04</b>	<b>4.38</b>	0.9178

Table 4: **Network-in-Network CIFAR-10 Results.** Accuracy, multiply-accumulate operations, and number of parameters for the baseline Network-in-Network model and more efficient versions created by the methods described in this paper.

This architecture uses  $5 \times 5$  filters in some layers. We found that we could replace all of these with  $3 \times 3$  filters, with comparable accuracy. As suggested by [Simonyan & Zisserman \(2014\)](#), stacked  $3 \times 3$  filters have the effective receptive field of larger filters with less computational complexity. In this **nin-c3** network, we replaced the first convolutional layer with one  $3 \times 3$  layer, and the second convolutional layer with two  $3 \times 3$  layers. This network is 26% faster than the standard NiN model, with only 54% of the model parameters. Using our low-rank filters in this network, we trained the **nin-c3-lr** network, which is of similar accuracy (91.8% v.s. 91.9%) but is approximately 54% of the original network’s computational complexity, with only 45% of the model parameters.

## 5 DISCUSSION

It is somewhat surprising that networks based on learning filters with less representational ability are able to do as well, or better, than CNNs with full  $k \times k$  filters on the task of image classification. However, a lot of interesting small-scale image structure is well-characterized by low-rank filters, *e.g.* edges and gradients. Our experiments training a separable (rank-1) model (‘vgg-gmp-sf’) on ILSVRC and MIT Places show surprisingly high accuracy on what are considered challenging problems – approx. 88% top-5 accuracy on ILSVRC – but not enough to obtain comparable accuracies to the models on which they are based.

Given that most discriminative filters learned for image classification appear to be low-rank, we instead structure our architectures with a set of basis filters in the way illustrated in Fig. 2d. This allows our networks to learn the most effective combinations of complex (*e.g.*  $k \times k$ ) and simple (*e.g.*  $1 \times k$ ,  $k \times 1$ ) filters. Furthermore, in restricting how many complex spatial filters may be learned, this architecture prevents over-fitting, and helps improve generalization. Even in our models where we do not use square  $k \times k$  filters, we obtain comparable accuracies to the baseline model, since the rank-2 cross-shaped filters effectively learned as a combination of  $3 \times 1$  and  $1 \times 3$  filters are capable of representing more complex local pixel relations than rank-1 filters.

## 6 CONCLUSION

This paper has presented a method to train convolutional neural networks from scratch using low-rank filters. This is made possible by a new way of initializing the networks weights which takes into consideration the presence of differently shape filters in composite layers. Validation on image classification in three popular datasets confirms similar or higher accuracy than state of the art, with much greater computational efficiency.

Recent advances in state-of-the-art accuracy with CNNs for image classification have come at the cost of increasingly large and computational complex models. We believe our results to show that learning computationally efficient models with fewer, more relevant parameters, can prevent over-fitting, increase generalization and thus also increase accuracy.

### FUTURE WORK

This paper has addressed the spatial extents of convolutional filters in CNNs, however the channel extents also exhibit some redundancy, as highlighted by [Jaderberg et al. \(2014\)](#), and exploited in the form of low-dimensional embeddings by [Lin et al. \(2013\)](#); [Szegedy et al. \(2014\)](#). We intend to further explore how our methods can be extended to learn and combine even smaller basis filters and filters with more diverse shapes.

## REFERENCES

- Bottou, Léon. Stochastic gradient descent tricks. In *Neural Networks: Tricks of the Trade*, pp. 421–436. Springer, 2012.
- Denil, Misha, Shakibi, Babak, Dinh, Laurent, de Freitas, Nando, et al. Predicting parameters in deep learning. In *Advances in Neural Information Processing Systems*, pp. 2148–2156, 2013.
- Fukushima, Kunihiro. Neocognitron: A Self-Organizing Neural Network Model for a Mechanism of Pattern Recognition Unaffected by Shift in Position. *Biological Cybernetics*, 36:193–202, 1980.
- Glorot, Xavier and Bengio, Yoshua. Understanding the difficulty of training deep feedforward neural networks. In *International conference on artificial intelligence and statistics*, pp. 249–256, 2010.
- Goodfellow, Ian, Warde-farley, David, Mirza, Mehdi, Courville, Aaron, and Bengio, Yoshua. Maxout Networks. In *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, pp. 1319–1327, 2013.
- Gupta, Suyog, Agrawal, Ankur, Gopalakrishnan, Kailash, and Narayanan, Pritish. Deep Learning with Limited Numerical Precision, February 2015.
- He, Kaiming, Zhang, Xiangyu, Ren, Shaoqing, and 0001, Jian Sun. Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification. *CoRR*, abs/1502.01852, 2015.
- Hochreiter, Sepp, Bengio, Yoshua, Frasconi, Paolo, and Schmidhuber, Jürgen. Gradient Flow in Recurrent Nets: the Difficulty of Learning Long-Term Dependencies, 2001.
- Jaderberg, Max, Vedaldi, Andrea, and Zisserman, Andrew. Speeding up Convolutional Neural Networks with Low Rank Expansions. *CoRR*, abs/1405.3866, 2014.
- Jia, Yangqing, Shelhamer, Evan, Donahue, Jeff, Karayev, Sergey, Long, Jonathan, Girshick, Ross, Guadarrama, Sergio, and Darrell, Trevor. Caffe: Convolutional Architecture for Fast Feature Embedding. *arXiv preprint arXiv:1408.5093*, 2014.
- Krizhevsky, Alex. Learning Multiple Layers of Features from Tiny Images. Technical report, University of Toronto, 2009.
- Krizhevsky, Alex, Sutskever, Ilya, and Hinton, Geoffrey E. ImageNet Classification with Deep Convolutional Neural Networks. In Pereira, F., Burges, C.J.C., Bottou, L., and Weinberger, K.Q. (eds.), *Advances in Neural Information Processing Systems 25*, pp. 1097–1105. Curran Associates, Inc., 2012.
- LeCun, Yann, Denker, John S, Solla, Sara A, Howard, Richard E, and Jackel, Lawrence D. Optimal brain damage. In *NIPs*, volume 89, 1989.
- LeCun, Yann, Bottou, Léon, Bengio, Yoshua, and Haffner, Patrick. Gradient-Based Learning Applied to Document Recognition. In *Proceedings of the IEEE*, volume 86, pp. 2278–2324, 1998.
- Lin, Min, Chen, Qiang, and Yan, Shuicheng. Network In Network. *CoRR*, abs/1312.4400, 2013.
- Mamalet, Franck and Garcia, Christophe. Simplifying convnets for fast learning. In *Artificial Neural Networks and Machine Learning–ICANN 2012*, pp. 58–65. Springer, 2012.
- Ren, Shaoqing, He, Kaiming, Girshick, Ross, Zhang, Xiangyu, and Sun, Jian. Object Detection Networks on Convolutional Feature Maps. *arXiv preprint arXiv:1504.06066*, 2015.
- Rigamonti, Roberto, Sironi, Amos, Lepetit, Vincent, and Fua, Pascal. Learning Separable Filters. In *CVPR*, pp. 2754–2761. IEEE, 2013.
- Simonyan, Karen and Zisserman, Andrew. Very Deep Convolutional Networks for Large-Scale Image Recognition. *CoRR*, abs/1409.1556, 2014.
- Szegedy, Christian, Liu, Wei, Jia, Yangqing, Sermanet, Pierre, Reed, Scott, Anguelov, Dragomir, Erhan, Dumitru, Vanhoucke, Vincent, and Rabinovich, Andrew. Going Deeper with Convolutions. *CoRR*, abs/1409.4842, 2014.
- Tompson, Jonathan, Goroshin, Ross, Jain, Arjun, LeCun, Yann, and Bregler, Christoph. Efficient Object Localization Using Convolutional Networks. June 2015.
- Vanhoucke, Vincent, Senior, Andrew, and Mao, Mark Z. Improving the speed of neural networks on CPUs. In *Proc. Deep Learning and Unsupervised Feature Learning NIPS Workshop*, volume 1, 2011.
- Xu, Li, Ren, Jimmy S, Liu, Ce, and Jia, Jiaya. Deep convolutional neural network for image deconvolution. In *Advances in Neural Information Processing Systems*, pp. 1790–1798, 2014.

# APPENDICES

## A INITIALIZING CNNs WITH MIXED-SHAPE LOW-RANK FILTERS

At the start of training, network weights are initialized at random using samples drawn from a Gaussian distribution with a standard deviation parameter specified separately for each layer. We found that the setting of these parameters was critical to the success of network training and difficult to get right, particularly because published parameter settings used elsewhere were not suitable for our new network architectures. With unsuitable weight initialization, training may fail due to *exploding gradients*, where back propagated gradients grow so large as to cause numeric overflow, or *vanishing gradients* where back propagated gradients grow so small that their effect is dwarfed by that of weight decay such that loss does not decrease during training (Hochreiter et al., 2001).

To determine the standard deviations to be used for weight initialization, we use an approach similar to that described by Glorot & Bengio (2010) (with the adaptation described by He et al. (2015) for layers followed by a ReLU). Their approach works by ensuring that the magnitudes of back-propagated gradients remain approximately the same throughout the network. Otherwise, if the gradients were inappropriately scaled by some factor (*e.g.*  $\beta$ ) then the final back-propagated signal would be scaled by a potentially much larger factor ( $\beta^L$  after  $L$  layers).

In what follows, we adopt notation similar to that of He et al. (2015), and follow their derivation of the appropriate standard deviation for weight initialization. However, we also generalize their approach to the initialization of ‘composite’ layers comprising several groups of filters of different spatial dimensions (see Fig. 5). This is one of the main contributions of this work.

**Forward Propagation.** The response of the  $l^{\text{th}}$  convolutional layer can be represented as

$$\mathbf{y}_l = \mathbf{W}_l \mathbf{x}_l + \mathbf{b}_l. \quad (2)$$

Here  $\mathbf{y}_l$  is a  $d \times 1$  vector representing a pixel in the output feature map and  $\mathbf{x}_l$  is a  $whc \times 1$  vector that represents a  $w \times h$  subregion of the  $c$ -channel input feature map.  $\mathbf{W}_l$  is the  $d \times n$  weight matrix, where  $d$  is the number of filters and  $n$  is the size of a filter, *i.e.*  $n = whc$  for a filter with spatial dimensions  $w \times h$  operating on an input feature map of  $c$  channels, and  $\mathbf{b}_l$  is the bias. Finally  $\mathbf{x}_l = f(\mathbf{y}_{l-1})$  is the output of the previous layer passed through an activation function  $f$  (*e.g.* the application of a ReLU to each element of  $\mathbf{y}_{l-1}$ ).

**Backward Propagation.** During back-propagation, the gradient of a convolutional layer is computed as

$$\Delta \mathbf{x}_l = \hat{\mathbf{W}}_l \Delta \mathbf{y}_l, \quad (3)$$

where  $\Delta \mathbf{x}_l$  and  $\Delta \mathbf{y}_l$  denote the derivatives of loss  $\mathcal{L}$  with respect to input and output pixels.  $\Delta \mathbf{x}_l$  is a  $c \times 1$  vector of gradients with respect to the channels of a single pixel in the input feature map and  $\Delta \mathbf{y}$  represents  $h \times w$  pixels in  $d$  channels of the output feature map.  $\hat{\mathbf{W}}_l$  is a  $c \times \hat{n}$  matrix where the filter weights are arranged in the right order for back-propagation, and  $\hat{n} = whd$ . Note that  $\hat{\mathbf{W}}_l$  can be simply reshaped from  $\mathbf{W}_l^T$ . Also note that the elements of  $\Delta \mathbf{y}_l$  correspond to pixels in the output image that had a forwards dependency on the input image pixel corresponding to  $\Delta \mathbf{x}$ . In back propagation, each element  $\Delta y_l$  of  $\Delta \mathbf{y}_l$  is related to an element  $\Delta x_{l+1}$  of some  $\Delta \mathbf{x}_{l+1}$  (*i.e.* a back-propagated gradient in the next layer) by the derivative of the activation function  $f$ :

$$\Delta y_l = f'(y_l) \Delta x_{l+1}, \quad (4)$$

where  $f'$  is the derivative of the activation function.

**Weight Initialization.** Now let  $\Delta y_l$ ,  $\Delta x_l$  and  $w_l$  be scalar random variables that describe the distribution of elements in  $\Delta \mathbf{y}_l$ ,  $\Delta \mathbf{x}_l$  and  $\hat{\mathbf{W}}_l$  respectively. Then, assuming  $f'(y_l)$  and  $\Delta x_{l+1}$  are independent,

$$E[\Delta y_l] = E[f'(y_l)]E[\Delta x_{l+1}]. \quad (5)$$

For the ReLU case,  $f'(y_l)$  is zero or one with equal probability. Like Glorot & Bengio (2010), we assume that  $w_l$  and  $\Delta y_l$  are independent of each other. Thus, equation 3 implies that  $\Delta x_l$  has zero mean for all  $l$ , when  $w_l$  is initialized by a distribution that is symmetric around zero. Thus we have

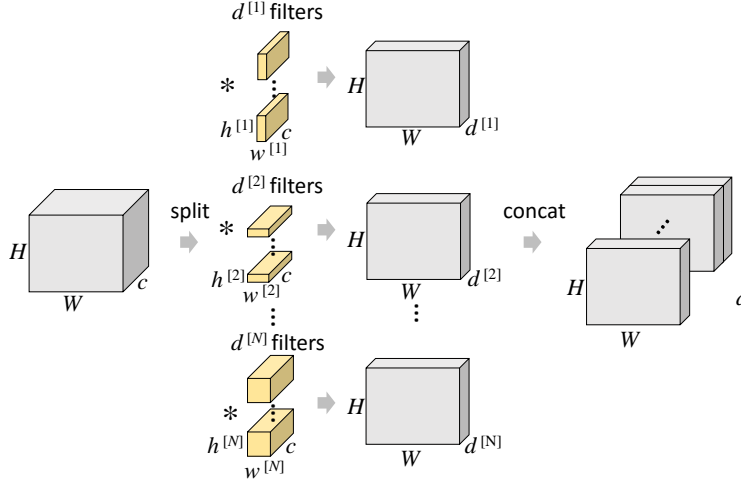


Figure 5: **A composite layer.** Composite layers convolve an input feature map with  $N$  groups of convolutional filters of several different spatial dimensions. Here the  $i^{\text{th}}$  group has  $d^{[i]}$  filters with spatial dimension  $w^{[i]} \times h^{[i]}$ . The outputs are concatenated to create a  $d$  channel output feature map. Composite layers require careful weight initialization to avoid vanishing/exploding gradients during training.

$E[\Delta y_l] = \frac{1}{2}E[\Delta x_{l+1}] = 0$  and also  $E[(\Delta y_l)^2] = \text{Var}[\Delta y_l] = \frac{1}{2}\text{Var}[\Delta x_{l+1}]$ . Now, since each element of  $\Delta \mathbf{x}_l$  is a summation of  $\hat{n}$  products of elements of  $\hat{\mathbf{W}}_l$  and elements of  $\Delta \mathbf{y}_l$ , we can compute the variance of the gradients in equation 3:

$$\begin{aligned} \text{Var}[\Delta x_l] &= \hat{n} \text{Var}[w_l] \text{Var}[\Delta y_l] \\ &= \frac{1}{2} \hat{n} \text{Var}[w_l] \text{Var}[\Delta x_{l+1}]. \end{aligned} \quad (6)$$

To avoid scaling the gradients in the convolutional layers (to avoid exploding or vanishing gradients), we set the ratio between these variances to 1:

$$\frac{1}{2} \hat{n} \text{Var}[w_l] = 1. \quad (7)$$

This leads to the result of He et al. (2015), in that a layer with  $\hat{n}_l$  connections followed by a ReLU activation function should be initialized with a zero-mean Gaussian distribution with standard deviation  $\sqrt{2/\hat{n}_l}$ .

**Weight Initialization in Composite Layers.** The initialization scheme described above assumes that the layer comprises filters of spatial dimension  $w \times h$ . Now we extend this scheme to composite convolutional layers containing  $N$  groups of filters of different spatial dimensions  $w^{[i]} \times h^{[i]}$  (where superscript  $[i]$  denotes the group index and with  $i \in \{1, \dots, N\}$ ). Now the layer response is the concatenation of the responses of each group of filters:

$$\mathbf{y}_l = \begin{bmatrix} \mathbf{W}_l^{[1]} \mathbf{x}_l^{[1]} \\ \mathbf{W}_l^{[2]} \mathbf{x}_l^{[2]} \\ \vdots \\ \mathbf{W}_l^{[N]} \mathbf{x}_l^{[N]} \end{bmatrix} + \mathbf{b}_l. \quad (8)$$

As before  $\mathbf{y}_l$  is a  $d \times 1$  vector representing the response at one pixel of the output feature map. Now each  $\mathbf{x}^{[i]}$  is a  $w^{[i]} h^{[i]} c \times 1$  vector that represents a different shaped  $w^{[i]} \times h^{[i]}$  sub-region of the input feature map. Each  $\mathbf{W}_l^{[i]}$  is the  $c_l^{[i]} \times \hat{n}^{[i]}$  weight matrix, where  $d$  is the number of filters and  $\hat{n}^{[i]}$  is the size of a filter, i.e.  $\hat{n}^{[i]} = w^{[i]} h^{[i]} c^{[i]}$  for a filter of spatial dimension  $w^{[i]} \times h^{[i]}$  operating on an input feature map of  $c_l = d_{l-1}$  channels.

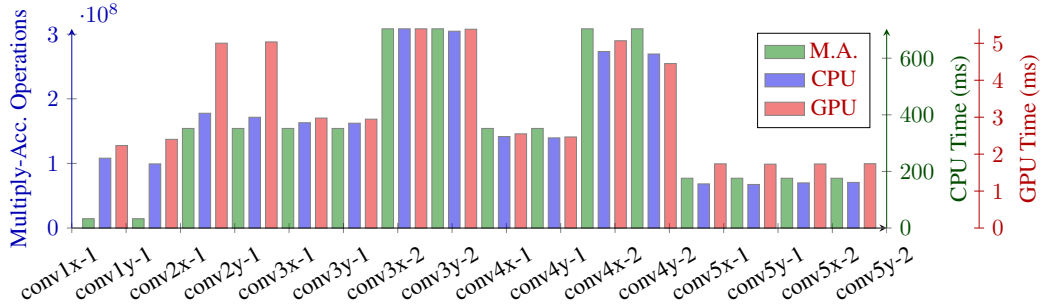


Figure 6: **Multiply-Accumulate Operations and Caffe CPU/GPU Timings.** For the forward pass of each convolutional layer in the ‘vgg-gmp-lr’ network. Caffe CPU and GPU timings were well correlated with multiply-accumulate operations for most layers.

During back propagation, the gradient of the composite convolutional layer is computed as a summation of the contributions from each group of filters:

$$\Delta \mathbf{x}_l = \hat{\mathbf{W}}_l^{[1]} \Delta \mathbf{y}_l^{[1]} + \hat{\mathbf{W}}_l^{[2]} \Delta \mathbf{y}_l^{[2]} + \dots + \hat{\mathbf{W}}_l^{[N]} \Delta \mathbf{y}_l^{[N]}, \quad (9)$$

where now  $\Delta \mathbf{y}^{[z]}$  represents  $w^{[z]} \times h^{[z]}$  pixels in  $d^{[z]}$  channels of the output feature map. Each  $\hat{\mathbf{W}}_l^{[z]}$  is a  $c_l \times \hat{n}^{[z]}$  matrix of weights arranged appropriately for back propagation. Again, note that each  $\hat{\mathbf{W}}_l^{[z]}$  can be simply reshaped from  $\mathbf{W}_l^{[z]}$ .

As before, each element of  $\Delta \mathbf{y}_l$  is a sum over  $\hat{n}$  products between elements of  $\hat{\mathbf{W}}_l^{[z]}$  and elements of  $\Delta \mathbf{y}_l^{[z]}$  and here  $\hat{n}$  is given by:

$$\hat{n} = \sum w^{[z]} h^{[z]} d^{[z]}. \quad (10)$$

In the case of a ReLU non-linearity, this leads to a zero-mean Gaussian distribution with standard deviation:

$$\sigma = \sqrt{\frac{2}{\sum w^{[z]} h^{[z]} d^{[z]}}}. \quad (11)$$

In conclusion, a composite layer of heterogeneously-shaped filter groups, where each filter group  $i$  has  $w^{[i]} h^{[i]} d^{[i]}$  outgoing connections should be initialized as if it is a single layer with  $\hat{n} = \sum w^{[i]} h^{[i]} d^{[i]}$ .

## B MULTIPLY-ACCUMULATE OPERATIONS AND CAFFE CPU/GPU TIMINGS.

We have characterized the computational complexity of a CNN using the number of multiply-accumulate operations required for a forward pass (which depends on the size of the filters in each convolutional layer as well as the input image size and stride), to give as close as possible to a hardware and implementation independent evaluation the computational complexity of our method. However, we have observed strong correlation between multiply-accumulate counts and runtime for both CPU and GPU implementations of the networks described here (as shown in Fig. 6). Note that the Caffe timings differ more for the initial convolutional layers where the input sizes are much smaller (3-channels), and BLAS is less efficient for the relatively small matrices being multiplied.

## C COMPARING WITH STATE OF THE ART NETWORKS FOR ILSVRC

Figures 7 and 8 compare published top-5 ILSVRC validation error v.s. multiply-accumulate operations and number of model parameters (respectively) for several state-of-the-art networks (Simonyan & Zisserman, 2014; Szegedy et al., 2014; He et al., 2015). The error rates for these networks are only reported as obtained with different combinations of computationally expensive training and

test time augmentation methods, including scale, photometric, ensembles (multi-model), and multi-view/dense oversampling. This can make it difficult to compare model architectures, especially with respect to computational requirements.

State-of-the-art networks, such as MSRA-C, VGG-19 and oversampled GoogLeNet are orders of magnitude larger in computational complexity than our networks. From Fig. 7, where the multiply-accumulate operations are plotted on a log scale, increasing the model size and/or computational complexity of test-time augmentation of CNNs appears to have diminishing returns for decreasing validation error. Our models *without* training or test time augmentation show comparable accuracy to networks such as VGG-13 *with* training and test time augmentation, while having far less computational complexity and model size. In particular, the ‘googlenet-lr’ model has a much smaller test-time model size than any network of comparable accuracy.

Network	Multiply-Acc. $\times 10^9$	Test M.A. w/ Aug. $\times 10^9$	Param. $\times 10^7$	Top-5 Acc.
msra-c	53.46	107.17	33.06	<b>0.943</b>
msra-b	23.22	46.54	18.33	0.937
msra-a	19.06	38.20	17.80	0.935
vgg-19	19.63	39.30	14.37	0.910
vgg-16 (D)	15.47	30.97	13.84	0.912
vgg-16 (C)	11.77	23.57	13.36	0.906
vgg-13	11.31	22.64	13.30	0.901
vgg-11	7.61	<b>15.24</b>	13.29	0.895
googlenet 10x	<b>1.59</b>	15.91	<b>1.34</b>	0.909
googlenet 144x	1.59	229.11	1.34	0.921

Table 5: **State of the Art Single Models with Extra Augmentation.** Top-5 ILSVRC validation accuracy, single view and augmented test-time multiply-accumulate (M.A.) count, and number of parameters for various state of the art models *with* various training and test-time augmentation methods. A multi-model ensemble of MSRA-C is the current state of the art network.

## D PLOTS OF RESULTS

Following are several plots of results, which for reasons of space consideration are not in the main section of the paper. These include the results for VGG-derived models on MIT Places (Fig. 9), GoogLeNet-derived models on ILSVRC (Fig. 10), and finally the results for Network-in-Network-derived models on CIFAR-10 (Fig. 11).

## E VGG-DERIVED MODEL TABLE

Table 6 shows the architectural details of the VGG-11 derived models used in §4.1.

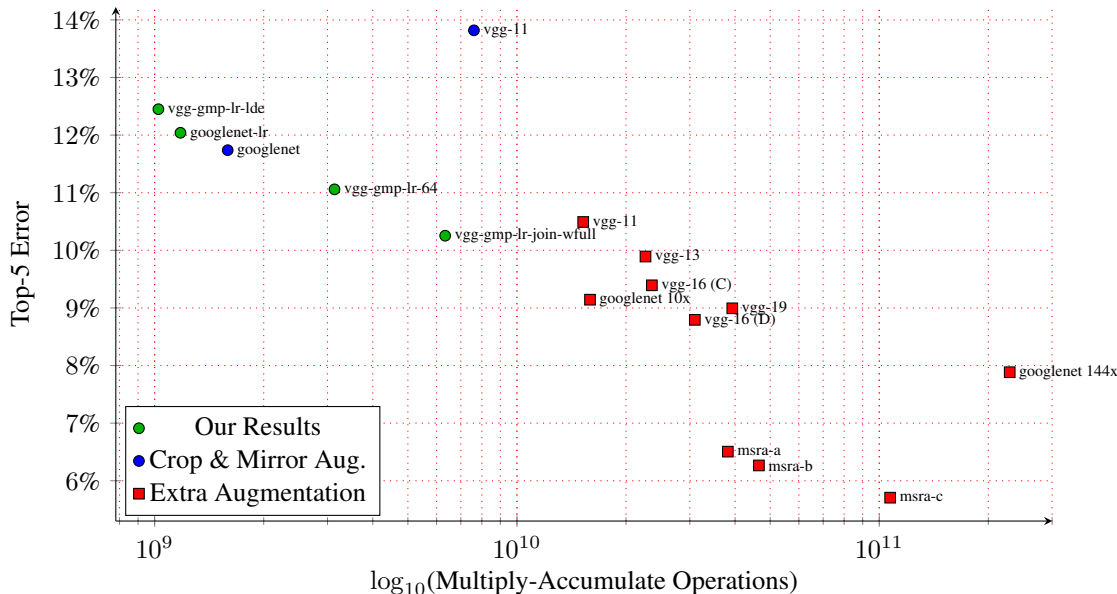


Figure 7: **Computational Complexity of Single State-of-the-Art ILSVRC Models.** Test-time multiply-accumulate operations v.s. top-5 error on state of the art networks using a *single* model. Note the difference in accuracy and computational complexity for VGG-11 model with/without extra augmentation. Our ‘vgg-gmp-lr-join-wfull’ model *without* extra augmentation is more accurate than VGG-11 *with* extra augmentation, and is much less computationally complex.

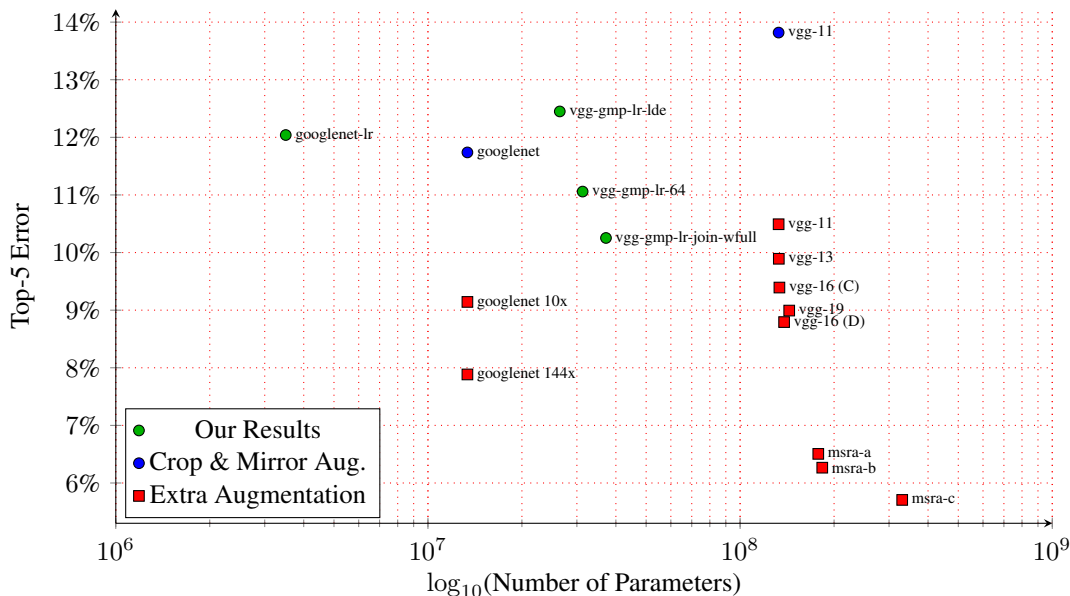


Figure 8: **Number of Parameters of State-of-the-Art ILSVRC Models.** Test time parameters v.s. top-5 error for state of the art models. The main factor in reduced model size is the use of global pooling or lack of fully-connected layers. Note that our ‘googlenet-lr’ model is almost an order of magnitude smaller than any other network of comparable accuracy.

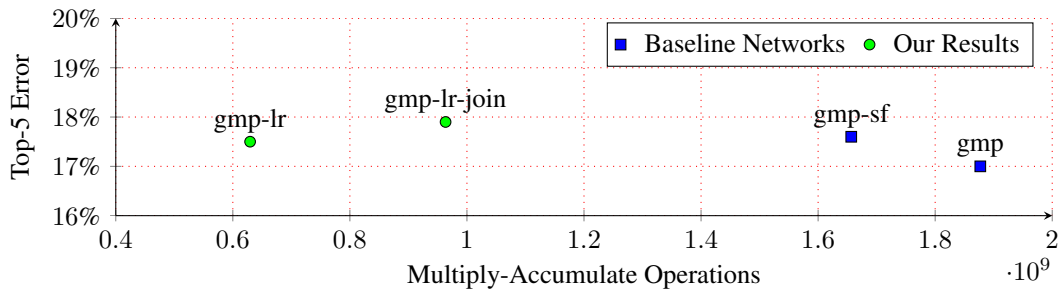


Figure 9: **MIT Places Results.** Multiply-accumulate operations v.s. top-5 error for VGG-derived models on MIT Places scene classification dataset.

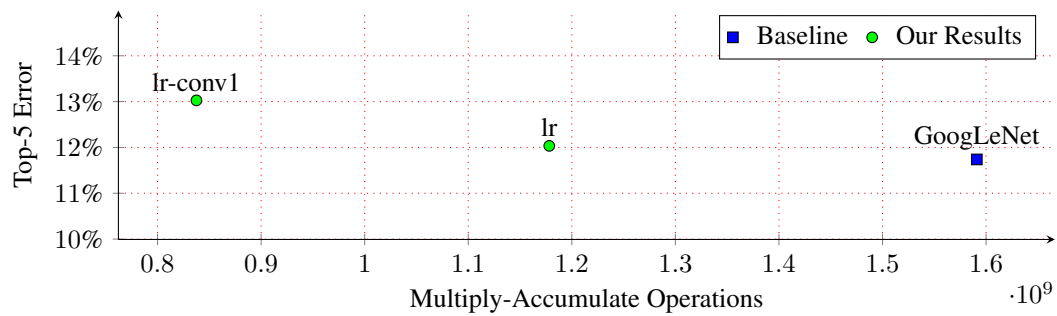


Figure 10: **GoogLeNet ILSVRC Results.** Multiply-accumulate operations v.s. top-5 error for GoogLeNet-derived models on ILSVRC object classification dataset.

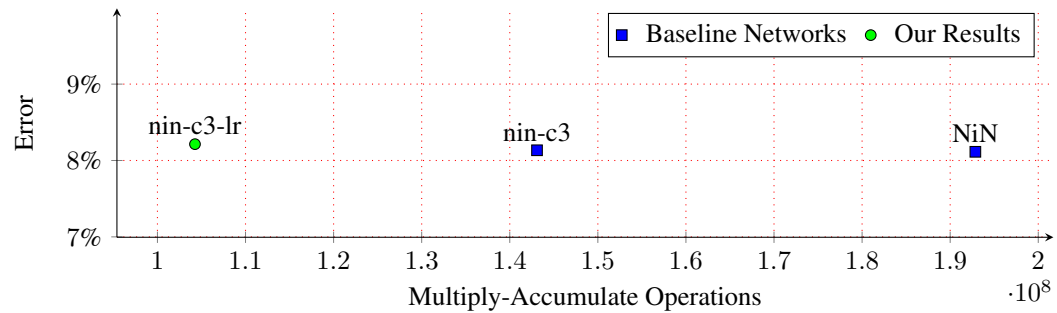


Figure 11: **Network-in-Network CIFAR-10 Results.** Multiply-accumulate operations v.s. error for Network-in-Network derived models on CIFAR-10 object classification dataset.



Layer	VGG-11	GMP	GMP-SF	GMP-LR	GMP-LR-2X	GMP-LR-JOIN	GMP-LR-LDE	GMP-LR-JOIN-WFULL	
conv1	3×3, 64	1×3, 64	3×1, 32    1×3, 32	3×1, 64    1×3, 64	3×1, 32    1×3, 32	3×1, 32    1×3, 32	3×1, 24    1×3, 24    3×3, 16		
		3×1, 64			1×1, 64			1×1, 32	1×1, 64
	ReLU								
	2×2 maxpool, /2								
conv2	3×3, 128	1×3, 128	3×1, 64    1×3, 64	3×1, 128    1×3, 128	3×1, 64    1×3, 64	3×1, 64    1×3, 64	3×1, 48    1×3, 48    3×3, 32		
		3×1, 128			1×1, 128			1×1, 64	1×1, 128
	ReLU								
	2×2 maxpool, /2								
conv3	3×3, 256	1×3, 256	3×1, 128    1×3, 128	3×1, 256    1×3, 256	3×1, 128    1×3, 128	3×1, 128    1×3, 128	3×1, 96    1×3, 96    3×3, 64		
		3×1, 256			1×1, 256			1×1, 128	1×1, 256
	ReLU								
	3×3, 256	1×3, 256	3×1, 128    1×3, 128	3×1, 256    1×3, 256	3×1, 128    1×3, 128	3×1, 128    1×3, 128	3×1, 96    1×3, 96    3×3, 64		
		3×1, 256			1×1, 256			1×1, 128	1×1, 256
	ReLU								
	2×2 maxpool, /2								
	conv4	3×3, 512	1×3, 512	3×1, 256    1×3, 256	3×1, 512    1×3, 512	3×1, 256    1×3, 256	3×1, 256    1×3, 256	3×1, 192    1×3, 192    3×3, 128	
3×1, 512			1×1, 512			1×1, 256			1×1, 512
ReLU									
3×3, 512		1×3, 512	3×1, 256    1×3, 256	3×1, 512    1×3, 512	3×1, 256    1×3, 256	3×1, 256    1×3, 256	3×1, 192    1×3, 192    3×3, 128		
		3×1, 512			1×1, 512			1×1, 256	1×1, 512
ReLU									
2×2 maxpool, /2									
conv5	3×3, 512	1×3, 512	3×1, 256    1×3, 256	3×1, 512    1×3, 512	3×1, 256    1×3, 256	3×1, 256    1×3, 256	3×1, 192    1×3, 192    3×3, 128		
		3×1, 512			1×1, 512			1×1, 256	1×1, 512
	ReLU								
	3×3, 512	1×3, 512	3×1, 256    1×3, 256	3×1, 512    1×3, 512	3×1, 256    1×3, 256	3×1, 256    1×3, 256	3×1, 192    1×3, 192    3×3, 128		
		3×1, 512			1×1, 512			1×1, 256	1×1, 512
	ReLU								
2×2 maxpool, /2	global maxpool								
fc6	7 <sup>2</sup> × 512 × 4096	512 × 4096							
ReLU									
fc7	4096 × 4096								
ReLU									
fc8	4096 × 1000								
softmax									

Table 6: **VGG Model Architectures.** Here “3×3, 32” denotes 32 3×3 filters, “/2” denotes stride 2, fc denotes fully-connected, and || denotes a concatenation within a composite layer.