

Exploiting Linear Structure Within Convolutional Networks for Efficient Evaluation

Anonymous Author(s)

Affiliation

Address

email

Abstract

We present techniques for speeding up the test-time evaluation of large convolutional networks, designed for object recognition tasks. These models deliver impressive accuracy, but each image evaluation requires millions of floating point operations, making their deployment on smartphones and Internet-scale clusters problematic. The computation is dominated by the convolution operations in the lower layers of the model. We exploit the redundancy present within the convolutional filters to derive approximations that significantly reduce the required computation. Using large state-of-the-art models, we demonstrate speedups by a factor of $2\times$, while keeping the accuracy within 1% of the original model.

1 Introduction

Large neural networks have recently demonstrated impressive performance on a range of speech and vision tasks. However, the size of these models can make their deployment at test time problematic. For example, mobile computing platforms are limited in their CPU speed, memory and battery life. At the other end of the spectrum, Internet-scale deployment of these models requires thousands of servers to process the 100's of millions of images per day. The electrical and cooling costs of these servers required is significant. Training large neural networks can take weeks, or even months. This hinders research and consequently there have been extensive efforts devoted to speeding up training procedure. However, there are relatively few efforts aimed at improving the *test-time* performance of the models.

We consider convolutional neural networks used for computer vision tasks, since they are large and widely used in commercial applications. These networks typically require a huge number of parameters ($\sim 10^7$ in [?]) to produce state-of-the-art results. This redundancy seems necessary in order to overcome a highly non-convex optimization [1, ?], but as a byproduct the resulting network wastes computing resources. In this paper we show that this redundancy can be exploited with linear compression techniques, resulting in significant speedups for the evaluation of *trained* large scale networks, with minimal compromise to performance. In particular, we concentrate in the lower convolutional layers, which typically dominate the evaluation cost.

We follow a relatively simple strategy: we start by compressing each convolutional layer by finding an appropriate low-rank approximation, and then we fine-tune the upper layers until the prediction performance is restored. We consider several elementary tensor decompositions based on singular value decompositions, as well as vector quantization to also exploit nonlinear approximation.

In summary, our main contributions are:

- We present a generic method to exploit the redundancy inherent in deep convolutional networks.

- We report experiments on state-of-the-art Imagenet CNNs, showing $2\times$ speedups on each convolutional layer with less than 1% drop in performance, both in GPU and CPU implementations.

Notation: Convolution weights can be described as a 4-dimensional tensor: $W \in \mathbb{R}^{C \times X \times Y \times F}$. C is the number of number of input channels, X and Y are the spatial dimensions of the kernel, and F is the target number of feature maps. It is common for the first convolutional layer to have a stride associated with the kernel. We denote the stride by Δ . Let $I \in \mathbb{R}^{C \times N \times M}$ denote an input signal where C is the number of input maps, and N and M are the spatial dimensions of the maps. The target value, $T = I * W$, of a generic convolutional layer, with $\Delta = 1$, for a particular output feature, f , and spatial location, (x, y) , is

$$T(f, x, y) = \sum_{c=1}^C \sum_{x'=1}^X \sum_{y'=1}^Y I(c, x - x', y - y') W(c, x', y', f)$$

Moreover, we define $W_C \in \mathbb{R}^{C \times (XYF)}$, and $W_F \in \mathbb{R}^{(CXY) \times F}$, and $W_S \in \mathbb{R}^{C \times (XY) \times F}$ to be the folded, with respect to different dimensions, versions of W .

If X is a tensor, $\|X\|$ denotes its operator norm, and $\|X\|_F$ denotes its Frobenius norm. If v is a vector, $\|v\|$ denotes its Euclidean norm.

2 Related Work

Vanhoucke *et al.* [2] explored the properties of CPUs to speed up execution. They present many solutions specific to Intel and AMD CPUs, however some of their techniques are general enough to be used for any type of processor. They describe how to align memory, and use SIMD operations (vectorized operations on CPU) to boost the efficiency of matrix multiplication. Additionally, they propose the linear quantization of the network weights and input. This involves representing weights as 8-bit integers (range $[-127, 128]$), rather than 32-bit floats. This approximation is similar in spirit to our approach, but differs in that it is applied to each weight element independently. By contrast, our approximation approach models the structure within each filter. Potentially, the two approaches could be used in conjunction.

The most expensive operations in convolutional networks are the convolutions in the first few layers. The complexity of this operation is linear in the area of the receptive field of the filters, which is relatively large for these layers. However, Mathieu *et al.* [3] have shown that convolution can be efficiently computed in Fourier domain, where it becomes element-wise multiplication (and there is no cost associated with size of receptive field). They report a forward-pass speed up of around $2\times$ for convolution layers in state-of-the-art models. Importantly, the FFT method can be used jointly with most of techniques presented in this paper.

The use of low-rank approximations in our approach is inspired by work of Denil *et al.* [1] who demonstrate the redundancies in neural network parameters. They show that the weights within a layer can be accurately predicted from a small (e.g. $\sim 5\%$) subset of them. This indicates that neural networks are heavily over-parametrized. All the methods presented here focus on exploiting the linear structure of this over-parametrization.

Finally, a recent preprint [?] also exploits low-rank decompositions of convolutional tensors to speed up the evaluation of convolutional neural networks, applied to scene text character recognition. This work was developed simultaneously with ours, and provides further evidence that such techniques can be applied to a variety of architectures and tasks.

3 Convolutional Tensor Compression

In this section we describe techniques to compress the 4 dimensional convolutional weight tensors into a representation that permits efficient computation. Section 3.1 describes how to construct good approximation criteria. Sections 3.2 and 3.4 describes techniques for low-rank tensor approximations and vector quantization.

3.1 Approximation Metric

The goal is to find an approximation \widetilde{W} of a convolutional tensor W which is cheaper to compute and such that the resulting network has similar prediction performance. Which metric should be used as approximation criteria?

A natural choice is to search for approximations such that $\|\widetilde{W} - W\|_F$ is small, since this criteria yields efficient compression schemes using elementary linear algebra, and also controls the operator norm of each linear convolutional layer. However, it assumes that all directions in the space of weights are equally affecting prediction performance. This metric can be improved while keeping the same efficient approximation algorithms.

We propose here a simple modification of the metric, of the form

$$\|W\|_\alpha^2 := \sum_p \alpha_p^2 W(p)^2, \quad (1)$$

where the sum runs over the tensor coordinates and $\alpha_p \geq 0$ are weights. Since (1) is a reweighted Euclidean metric, we can simply compute $W' = \alpha \cdot W$, where \cdot denotes element-wise multiplication, then compute the approximation \widetilde{W}' on W' using the standard L_2 norm, and finally output $\widetilde{W} = \alpha^{-1} \cdot \widetilde{W}'$. The natural question is then how to choose the weights α_p . For that purpose, we seek to emphasize coordinates more prone to produce prediction errors over coordinates whose effect is less harmful for the overall system.

We can obtain such measurements as follows. Let $\Theta = \{W_1, \dots, W_S\}$ denote the set of all parameters of the S -layer network, and let $U(I; \Theta)$ denote the output after the softmax layer of input image I . We consider a given input training set (I_1, \dots, I_N) with known labels (y_1, \dots, y_N) . For each pair (I_n, y_n) , we compute the forward propagation pass $U(I_n, \Theta)$, and define as $\{\beta_n\}$ the indices of the h largest values of $U(I_n, \Theta)$ different from y_n . Then, for a given layer s , we compute

$$d_{n,l,s} = \nabla_{W_s} (U(I_n, \Theta) - \delta(i - l)) , \quad n \leq N, l \in \{\beta_n\}, s \leq S, \quad (2)$$

where $\delta(i - l)$ is the dirac distribution centered at l . In other words, for each input we back-propagate the difference between the current prediction and the h ‘‘most dangerous’’ mistakes. The resulting weights α_p are then obtained by computing the average energy in the tensors $d_{n,l,s}$: $\alpha_p = \left(\sum_{n,l} d_{n,l,s}(p)^2 \right)^{1/2}$.

An even better metric can be obtained by considering the Mahalanobis distance defined from the covariance of d : $\|W\|_{maha}^2 = w \Sigma^{-1} w^T$, where w is the vector containing all the coordinates of W , and Σ is the covariance of $(d_{n,l,s})_{n,l}$. However, we do not report results using this metric, since it requires inverting a matrix of size equal to the number of parameters, which can be prohibitively expensive in large networks.

One can view the Frobenius norm of W as $\|W\|_F^2 = \mathbb{E}_{x \sim \mathcal{N}(0, I)} \|Wx\|_F^2$. Another alternative, which has also been considered in [?], is to replace the isotropic covariance assumption by the empirical covariance of the input of the layer. If $W \in \mathbb{R}^{C \times X \times Y \times F}$ is a convolutional layer, and $\widehat{\Sigma} \in \mathbb{R}^{CXY \times CXY}$ is the empirical estimate of the input data covariance, it can be efficiently computed as

$$\|W\|_{data}^2 = \|\widehat{\Sigma}^{1/2} W_f\|_F^2, \quad (3)$$

where W_f is the matrix obtained by folding the first three dimensions of W .

3.2 Low-rank Tensor Approximations

A particularly simple strategy to exploit the redundancy present in trained convolutional network weights is to linearly compress the tensors, which amounts to finding low-rank approximations.

Matrices are 2-tensors which can be linearly compressed using the Singular Value Decomposition. If $W \in \mathbb{R}^{m \times k}$ is a real matrix, it is defined as

$$W = U S V^T, \text{ where } U \in \mathbb{R}^{m \times m}, S \in \mathbb{R}^{m \times k}, V \in \mathbb{R}^{k \times k}.$$

S is a diagonal matrix with the singular values s_1, \dots, s_k on the diagonal, and U, V are orthogonal matrices. If the singular values of W decay rapidly, W can be well approximated by keeping only the t largest entries of S , resulting in the approximation

$$\tilde{W} = \tilde{U} \tilde{S} \tilde{V}^\top, \text{ where } \tilde{U} \in \mathbb{R}^{m \times t}, \tilde{S} \in \mathbb{R}^{t \times t}, \tilde{V} \in \mathbb{R}^{t \times k} \quad (4)$$

The approximation error $\|I\tilde{W} - IW\|_F$ satisfies

$$\|I\tilde{W} - IW\|_F \leq s_{t+1} \|I\|_F, \quad (5)$$

and thus is controlled by the decay along the diagonal of S . Now the computation $I\tilde{W}$ can be done in $O(nmt + nt^2 + ntk)$, which, for sufficiently small t is significantly smaller than $O(nmk)$.

The SVD can be used to approximate a tensor $W \in \mathbb{R}^{m \times n \times k}$ by first folding all but two dimensions together to convert it into a 2-tensor, and then considering the SVD of W_f .

Alternatively, the linear approximation of matrices can be easily extended to higher order tensors [?]. Let $v \otimes l$ denote the outer product of two vectors v and l . For a 3-tensor, $W_S \in \mathbb{R}^{C \times (XY) \times F}$, we can construct a rank 1 approximation by finding a decomposition that minimizes

$$\|W_S - \alpha \otimes \beta \otimes \gamma\|_F, \quad (6)$$

where $\alpha \in \mathbb{R}^C, \beta \in \mathbb{R}^{XY}, \gamma \in \mathbb{R}^F$ and $\|X\|_F$ denotes the Frobenius norm. Problem (6) is solved efficiently by performing alternate least squares on α, β and γ respectively, although more efficient algorithms can also be considered [?].

This easily extends to a rank K approximation using a greedy algorithm: Given a tensor M , we compute (α, β, γ) using (6), and we update $W_S \leftarrow W_S - \alpha \otimes \beta \otimes \gamma$. Repeating this operation K times results in

$$\tilde{W}_S = \sum_{k=1}^K \alpha_k \otimes \beta_k \otimes \gamma_k. \quad (7)$$

where $\alpha \in \mathbb{R}^C, \beta \in \mathbb{R}^{XY}, \gamma \in \mathbb{R}^F$. The approximations (6) and (7) are extended to q -tensors by adding more terms in the separable approximations. As opposed to the SVD for matrices, the rank-1 tensors are not orthogonal in general.

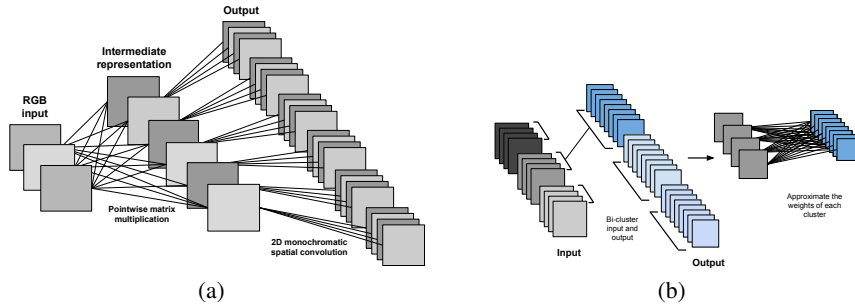


Figure 1: A visualization of monochromatic and biclustering approximation structures. (a): The monochromatic approximation, used for the first layer. Input color channels are projected by a set of intermediate color channels. After this transformation, output features need only to look at one intermediate color channel. This sparsity structure is what makes the speedups feasible. (b): The biclustering approximation, used for higher convolution layers. Input and output features are clustered into equal sized groups. The weight tensor corresponding to each pair of input and output clusters is then approximated.

3.3 Monochromatic Convolution Approximation

Let $W \in \mathbb{R}^{C \times X \times Y \times F}$ denote the weights of the first convolutional layer of a trained network. The number of input channels, C , corresponds to a different color component (either RGB or YUV). We found that the color components of trained convolutional neural networks we considered (see section 4) tend to have low dimensional structure. In particular, the weights can be well approximated

by projecting the color dimension down to a 1D subspace, i.e. a single color channel. The low-dimensional structure of the weights is apparent in 4.1.1. This figure shows the original first layer convolutional weights of a trained network and the weights after the color dimension has been projected into 1D lines.

The monochromatic approximation exploits this structure and is computed as follows. First, for every output feature, f , we consider the matrix $W_f \in \mathbb{R}^{C \times (X \cdot Y)}$, where the spatial dimensions of the filter corresponding to the output feature have been combined, and find the singular value decomposition, $W_f = U_f S_f V_f^\top$, where $U_f \in \mathbb{R}^{C \times C}$, $S_f \in \mathbb{R}^{C \times XY}$, and $V_f \in \mathbb{R}^{XY \times XY}$. We then take the rank 1 approximation of W_f :

$$\tilde{W}_f = \tilde{u}_f \tilde{s}_f \tilde{v}_f^\top, \quad (8)$$

where $\tilde{u}_f \in \mathbb{R}^{C \times 1}$, $\tilde{s}_f \in \mathbb{R}$, $\tilde{v}_f \in \mathbb{R}^{1 \times XY}$. This approximation corresponds to using a single color channel for each output feature. We can further exploit the regularity in the weights by sharing the color component basis between different output features. We do this by clustering the F left singular vectors, \tilde{u}_f , of each output feature f into C' clusters, where C' is much smaller than F . We constrain the clusters to be of equal size as discussed in section ???. Then, for each of the $\frac{F}{C'}$ output features f that is assigned to cluster c_f , we can approximate W_f with

$$\tilde{W}_f = u_{c_f} \tilde{s}_f \tilde{v}_f^\top \quad (9)$$

where $u_{c_f} \in \mathbb{R}^{C \times 1}$ is the cluster center for cluster c_f and \tilde{s}_f and \tilde{v}_f as as before.

This monochromatic approximation is illustrated in the left panel of Figure 1(a)). Table 1 shows the number of operations required for the standard and monochromatic versions.

3.4 Tensor Vector Quantization

Another source of redundancy within the 4-D weight tensors that can be exploited is the similarity between different filters. It can be efficiently captured by clustering the filters, such that each cluster can be accurately approximated by a low-rank factorization. For the filters in the upper layers, in practice this approach is more efficient than attempting to fit a low rank approximation to all filters.

On structured tensors such as those appearing in convolutional networks, we found it to be most efficient to cluster over both input and output feature channels using a standard biclustering scheme. The input clusters and output clusters are determined independently of one another. We start by clustering the rows of $W_C \in \mathbb{R}^{C \times (XYF)}$, which results in clusters C_1, \dots, C_a . Then we cluster the columns of $W_F \in \mathbb{R}^{(CXY) \times F}$. This procedure returns clusters F_1, \dots, F_b . These two operations break the original weight tensor W into ab sub-tensors $\{W_{C_i, F_j}\}_{i=1, \dots, a, j=1, \dots, b}$ as shown in Figure 1(b)). Each sub-tensor contains similar elements, thus should be easier to fit with a low-rank approximation.

To obtain a significant speedup of the test-time computation we must efficiently exploit the parallelism inherent in CPU and GPU architectures. The speedup obtainable is lower-bounded by the largest amount of computation assigned to a single thread. This implies that best strategy is to assign to every thread the same amount of work, or in other words have the clusters with the same number of filters. We therefore perform the biclustering operations (or clustering for monochromatic filters ??) using a modified version of the k-means algorithms which balances the cluster count at each iteration. It is implemented with the Floyd algorithm, by modifying the Euclidean distance with a subspace projection distance.

The number of input and output feature planes is large for all layers beyond the first. As described in Section ?? we cluster input and output features and then,

We then find low-rank approximations of each input-output subtensor W_{C_i, F_j} . For that purpose, we explore two different approximations previously described: (i) SVD by folding the tensor into a matrix, and (ii) outer product decomposition.

Using these approximations on the convolutional weights, the target output can be computed with significantly fewer operations. The number of operations required is a function of both the number of input clusters, G , and output clusters H . Moreover, let K denote the rank of the approximated tensors for the outer product decomposition. For the SVD decomposition, let K_1 denote the input

Approximation technique	Number of operations
No approximation	$XYCFNM\Delta^{-2}$
Monochromatic	$C'CNM + XYFNM\Delta^{-2}$
Biclustering + outer product decomposition	$GHK(NM\frac{C}{G} + XYNM\Delta^{-2} + \frac{F}{H}NM\Delta^{-2})$
Biclustering + SVD	$GHNM(\frac{C}{G}K_1 + K_1XYK_2\Delta^{-2} + K_2\frac{F}{H})$

Table 1: Number of operations required for various approximation methods.

mapping rank, and let K_2 denote the output mapping rank. The two different approximation methods have the following number of operations:

4 Experiments

We use two different convolutional architectures, each trained on the ImageNet 2012 dataset [4]: PierreNet [5], and MattNet [6]. We evaluate these networks on both CPU and GPU platforms.

We now present results showing the performance of the approximations described in section ?? in terms of prediction accuracy, speedup gains and reduction in memory overhead.

4.1 Speedup

Table 2 shows the time breakdown of forward propagation for each layer in MattNet. The majority of time is spent in the first and second layer convolution operations, so restrict our attention to these layers, although our approximations could easily be applied to convolutions in upper layers also.

We implemented several different approximation routines for the CPU and as well as the GPU in an effort to achieve empirical speedups. Both the baseline and approximation CPU code is implemented in C++ using Eigen3 library [7] compiled with Intel MKL. We also use Intel’s implementation of openmp, and multithreading. The baseline gives comparable performance to highly optimized MATLAB convolution routines and all of our CPU speedup results are computed relative to this. We used Alex Krizhevsky’s CUDA convolution routines ¹ as a baseline for GPU comparisons. The approximation versions are written in CUDA. All GPU code was run on a standard nVidia Titan card.

We have found that in practice it is often difficult to achieve speedups close to the theoretical gains based on the number of arithmetic operations. Moreover, different computer architectures and convolutional network architectures afford different optimization strategies making most implementations highly specific. However, regardless of implementation details, all of the approximations we present reduce both the number of operations and number of weights required to compute the output by at least a factor of two. While we present both theoretical and empirical speedup results, the empirical aspects serves largely as a proof of concept. Our aim is to illustrate that the reduction in floating point operations and memory accesses that are afforded by the approximations can in fact translate into speedups in practice. Our speedup results on both CPU and GPU are promising, especially given that our baselines for comparison comprise of highly optimized code. However, we believe that with extensive engineering efforts, speedups can be much closer to theoretically optimal levels.

4.1.1 First Layer

The first convolutional layer in MattNet has 3 input channels, 96 output channels and 7x7 filters. We approximated the weights in this layer using the monochromatic approximation described in section ?? . The monochromatic approximation works well if the color components span a small number of one dimensional subspaces. Figure 4.1.1 illustrates that by plotting the RGB components of the original filters and the filters after projecting onto one dimensional subspaces. Figure 4.1.1 gives an alternate presentation of this approximation, showing the original first layer filters beside their monochromatic approximations. Interestingly, we notice that the approximated filters often appear to be cleaned up versions of the original filters. This leads up to believe that the approximation techniques presented here might be viable strategies of cleaning up or denoising weights after training, potentially to improve generalization performance.

We evaluated MattNet on 8K validation images from the ImageNet12 dataset using various monochromatic approximations for the first layer weights. The only parameter in the approxima-

¹<https://code.google.com/p/cuda-convnet/>

Layer	Time per batch (sec)	Fraction
Conv1	2.8317 \pm 0.1030	21.97%
MaxPool	0.1059 \pm 0.0154	0.82%
LRNormal	0.1918 \pm 0.0162	1.49%
Conv2	4.2626 \pm 0.0740	33.07%
MaxPool	0.0705 \pm 0.0029	0.55%
LRNormal	0.0772 \pm 0.0027	0.60%
Conv3	1.8689 \pm 0.0577	14.50%
MaxPool	0.0532 \pm 0.0018	0.41%
Conv4	1.5261 \pm 0.0386	11.84%
Conv5	1.4222 \pm 0.0416	11.03%
MaxPool	0.0102 \pm 0.0006	0.08%
FC	0.3777 \pm 0.0233	2.93%
FC	0.0709 \pm 0.0038	0.55%
FC	0.0168 \pm 0.0018	0.13%
Softmax	0.0028 \pm 0.0015	0.02%
Total	12.8885	

Layer	Time per batch (sec)	Fraction
Conv1	0.0604 \pm 0.0112	5.14%
MaxPool	0.0072 \pm 0.0040	0.61%
LRNormal	0.0041 \pm 0.0043	0.35%
Conv2	0.4663 \pm 0.0072	39.68%
MaxPool	0.0032 \pm 0.0000	0.27%
LRNormal	0.0015 \pm 0.0003	0.13%
Conv3	0.2219 \pm 0.0014	18.88%
MaxPool	0.0016 \pm 0.0000	0.14%
Conv4	0.1991 \pm 0.0001	16.94%
Conv5	0.1958 \pm 0.0002	16.66%
MaxPool	0.0005 \pm 0.0001	0.04%
FC	0.0077 \pm 0.0013	0.66%
FC	0.0017 \pm 0.0001	0.14%
FC	0.0007 \pm 0.0002	0.06%
Softmax	0.0038 \pm 0.0098	0.32%
Total	1.1752	

Table 2: Evaluation time in seconds per layer of MattNet on CPU (left) and GPU (right) with batch size of 128. Results are averaged over 8 runs.

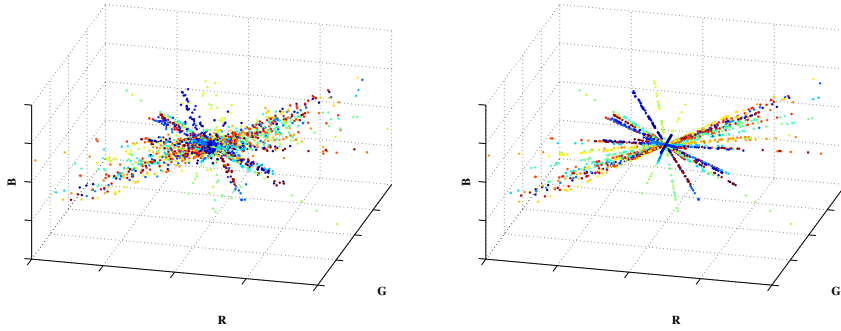


Figure 2: Visualization of the 1st layer filters in MattNet. Each component of the 96 7x7 filters is plotted in RGB space. Points are colored based on the output filter they belong to. Hence, there are 96 colors and 7^2 points of each color. (Left) Shows the original filters and (Right) shows the filters after the monochromatic approximation, where each filter has been projected down to a line in colorspace.

tion is C' , the number of color channels used for the intermediate representation. As expected, the network performance begins to degrade as C' decreases. Table 3 shows how the classification performance degrades for increasing C' .

Theoretically achievable speedups can be measured in terms of the number of floating point operations required to compute the output of the convolution. For the monochromatic approximation this theoretical speedup with respect to MattNet is given in table 3. The majority of the operations result from the convolution part of the computation. In comparison, the number of operations required for the color transformation is negligible. Thus, the theoretically achievable speedup decreases only slightly as the number of color components used is increased.

In practice, we found it difficult to optimize the monochromatic convolution routines as for large C' . Less work can be shared amongst output filters when C' is large making it challenging to parallelize. Figure 4 shows the empirical speedups we achieved on CPU (left) and GPU (right) and the corresponding network performance for various numbers of colors used in the monochromatic approximation. Our CPU and GPU implementations achieve empirical speedups of $\sim 1.8\times$ relative to the baseline with a drop in classification performance of less than 1.2%.

4.1.2 Second Layer

The second convolutional layer in MattNet has 96 input channels, 256 output channels and 5x5 filters. We approximated the weights using the biclustering operation, followed by the outer product decomposition approximation described in section ???. We evaluated the original and approximate convolution operation in this layer using 8K validation images from the ImageNet12 dataset. We explored various configurations of the approximations by varying the number of input clusters G , the number of output clusters H and the rank of the approximation.

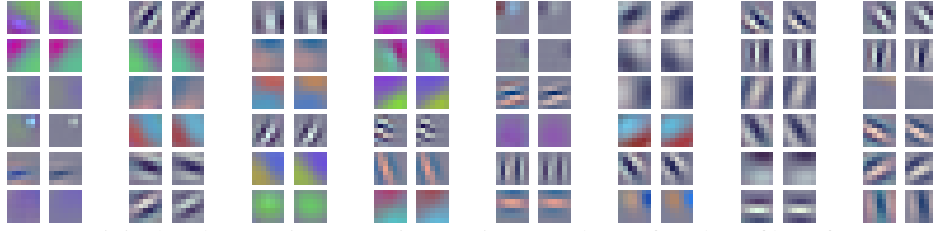


Figure 3: Original and approximate versions (using 12 colors) of 1st later filters from MattNet.

Number of colors	Increase in test error	Theoretical speedup
6	10.1929%	2.95×
8	4.0528%	2.94×
12	2.1973%	2.91×
16	1.1475%	2.88×
24	1.7212%	2.82×
48	0.5249%	2.66×
96	0.5738%	2.39×

Table 3: Performance of MattNet when first layer weights are replaced with monochromatic approximation and the corresponding theoretical speedup. Classification error on 8K validation images tends to increase as the approximation becomes harsher (i.e. fewer colors are used). Theoretical speedups vary only slightly as the number of colors used increases since the color transformation contributes relatively little to the total number of operations.

We can measure the theoretically achievable speedups for a particular approximation in terms of the number of floating point operations required to compute the output of the convolution. Figure 5 plots the theoretically achievable speedups against the drop in classification performance for various configurations. For a given setting of input and output clusters numbers, the performance tends to degrade as the rank is decreased.

In practice we found it difficult to achieve speedups close to theoretically optimal levels. This is unsurprising given the amount of engineering effort that has gone into the baseline routines. However, we achieved promising results and present speedups of $\sim 1.7\times$ relative to the baseline with less than a 1% drop in performance. Figure 6 shows our empirical speedups on CPU (left) and GPU (right) and the corresponding network performance for various approximation configurations. For the CPU implementation we used the biclustering with SVD decomposition approximation. For the GPU implementation we used the biclustering with outer product decomposition approximation. While we were unable to attain speedups on the order of those theoretically achievable, we are confident that further engineering efforts should decrease the margin.

4.2 Reduction in memory overhead

In many commercial applications memory conservation and storage are a central concern. This mainly applies to embedded systems (e.g. smartphones), where available memory is limited, and users are reluctant to download large files. In these cases, being able to compress the neural network is crucial for the viability of the product.

In addition to requiring fewer operations, our approximations require significantly fewer parameters when compared to the original model. Table 4 shows the number of parameters for various approximation methods as a function of hyperparameters for the approximation techniques. Table 5 shows the empirical reduction of parameters for MattNet and the corresponding network performance.

5 Discussion

In this paper we have presented techniques that can speed up the bottleneck convolution operations in the first layers of a convolutional network by a factor 1.6 – 2, with a moderate loss of performance ($\sim 1\%$). The empirical speedups achieved are still some way short of the theoretical gains, thus further improvements might be expected with further engineering effort. Moreover, the techniques are orthogonal to other approaches for efficient evaluation, such as quantization or working in the Fourier domain. Hence, they can potentially be used together to obtain further gains.

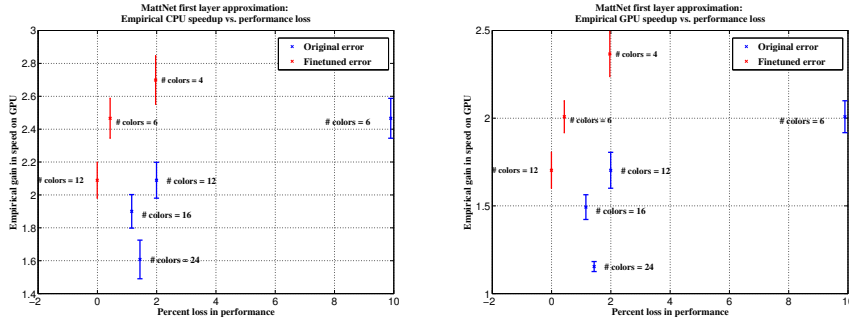


Figure 4: Empirical speedups on **(Left)** CPU and **(Right)** GPU for the first layer of MattNet.

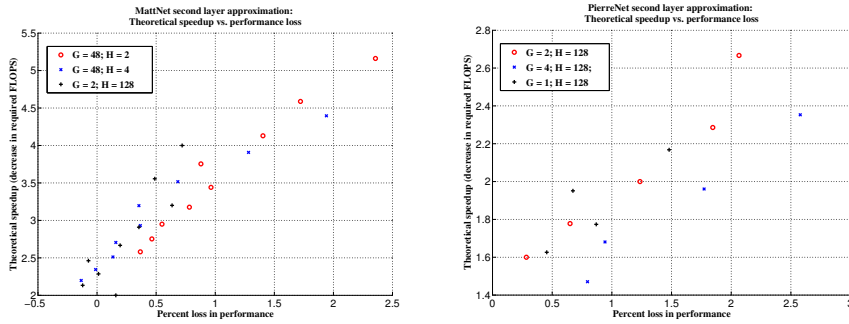


Figure 5: Theoretically achievable speedups with various biclustering approximations applied to **(Left)** MattNet and **(Right)** PierreNet.

Given the widespread use of neural networks in large scale industrial settings, the speed gains we demonstrate have significant economic value. Companies such as Google or Facebook process 10^8 images/day, requiring many thousands of machines. Our speed gains mean that the number of machines dedicated to running convolutional networks could be roughly halved, so saving millions of dollars, as well as a significant reduction in environmental impact.

We also show that our methods reduce the memory footprint of weights in the first two layers by factor of $2 - 3\times$. When applied to the whole network, this would translate into a significant savings, which would facilitate mobile deployment of convolutional networks.

Our approximations are applied to fully trained networks. The small performance drops that result could well be mitigated by further training the network after applying the approximation. By alternating these two steps, we could potentially achieve further gains in speed-up for a modest performance drop.

Another aspect of our technique is regularization. It seems that approximated filters look cleaner, and that sporadically we get better test error (e.g. bottom left of Figure 5(left)). We would like to experiment with low-rank projections during training as a regularization technique. Effectively it decreases number of learnable parameters, so it might improve generalization, a major issue with large convolutional networks.

Acknowledgments

References

- [1] Denil, M., Shakibi, B., Dinh, L., Ranzato, M., de Freitas, N.: Predicting parameters in deep learning. arXiv preprint arXiv:1306.0543 (2013)
- [2] Vanhoucke, V., Senior, A., Mao, M.Z.: Improving the speed of neural networks on cpus. In: Proc. Deep Learning and Unsupervised Feature Learning NIPS Workshop. (2011)

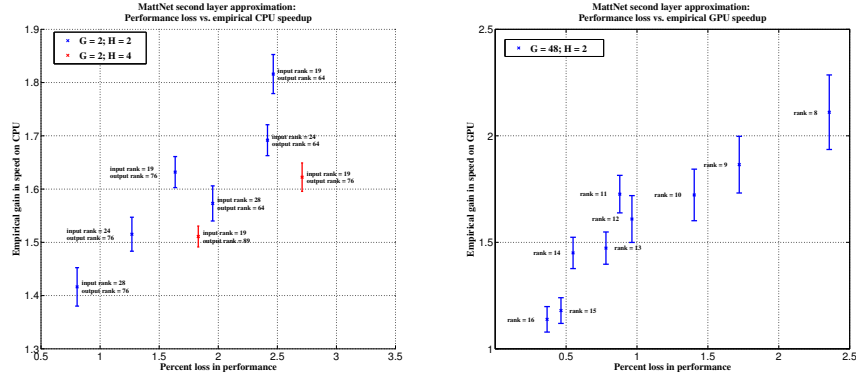


Figure 6: Empirical speedups on (Left) CPU and (Right) GPU for the second layer of MattNet.

Approximation method	Number of parameters
No approximation	$C \cdot X \cdot Y \cdot F$
Monochromatic	$C \cdot C' + X \cdot Y \cdot F$
Biclustering + outer product decomposition	$G \cdot H \cdot K \left(\frac{C}{G} + X \cdot Y + \frac{F}{H} \right)$
Biclustering + singular value decomposition	$G \cdot H \cdot \left(\frac{C}{G} \cdot K_1 + K_1 \cdot X \cdot Y \cdot K_2 + K_2 \cdot \frac{F}{H} \right)$

Table 4: Number of parameters expressed as a function of hyperparameters for various approximation methods.

- [3] Mathieu, M., Henaff, M., LeCun, Y.: Fast training of convolutional networks through ffts. arXiv preprint arXiv:1312.5851 (2013)
- [4] Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: ImageNet: A Large-Scale Hierarchical Image Database. In: CVPR09. (2009)
- [5] Sermanet, P., Eigen, D., Zhang, X., Mathieu, M., Fergus, R., LeCun, Y.: Overfeat: Integrated recognition, localization and detection using convolutional networks. arXiv preprint arXiv:1312.6229 (2013)
- [6] Zeiler, M.D., Fergus, R.: Visualizing and understanding convolutional neural networks. arXiv preprint arXiv:1311.2901 (2013)
- [7] Guennebaud, G., Jacob, B., et al.: Eigen v3. <http://eigen.tuxfamily.org> (2010)
- [8] Zeiler, M.D., Taylor, G.W., Fergus, R.: Adaptive deconvolutional networks for mid and high level feature learning. In: Computer Vision (ICCV), 2011 IEEE International Conference on, IEEE (2011) 1818–1825
- [9] Le, Q.V., Ngiam, J., Chen, Z., Chia, D., Koh, P.W., Ng, A.Y.: Tiled convolutional neural networks. In: Advances in Neural Information Processing Systems. (2010)
- [10] Le, Q.V., Ranzato, M., Monga, R., Devin, M., Chen, K., Corrado, G.S., Dean, J., Ng, A.Y.: Building high-level features using large scale unsupervised learning. arXiv preprint arXiv:1112.6209 (2011)
- [11] Lowe, D.G.: Object recognition from local scale-invariant features. In: Computer vision, 1999. The proceedings of the seventh IEEE international conference on. Volume 2., Ieee (1999) 1150–1157
- [12] Hinton, G.E., Srivastava, N., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.R.: Improving neural networks by preventing co-adaptation of feature detectors. arXiv preprint arXiv:1207.0580 (2012)
- [13] Krizhevsky, A., Sutskever, I., Hinton, G.: Imagenet classification with deep convolutional neural networks. In: Advances in Neural Information Processing Systems 25. (2012) 1106–1114

Approximation method	Approximation hyperparameters	Reduction in weights	Increase in test error
Layer 1: Monochromatic	$C' = 48$	$2.9109 \times$	0.5249%
Layer 2: Biclustering	$G = 48$	$3.7537 \times$	0.8789%
+outer product decomposition	$H = 2$		
	$K = 11$		
Layer 2: Biclustering	$G = 2$	$2.0386 \times$	0.8057%
+ singular value decomposition	$H = 2$		
	$K_1 = 28$		
	$K_2 = 76$		

Table 5: Empirical reduction in parameters for first two layers of MattNet with corresponding network performance.