
Deep Learning Approximation: Zero-Shot Neural Network Speedup

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Neural networks offer high-accuracy solutions to a range of problems, but are
2 computationally costly to run in production systems. We propose a technique
3 called Deep Learning Approximation to take an already-trained neural network
4 model and build a faster (and almost equally accurate) network by manipulating
5 the network structure and coefficients without requiring re-training or access to
6 the training data. Speedup is achieved by applying a sequential series of indepen-
7 dent optimizations that reduce the floating-point operations (FLOPs) required to
8 perform a forward pass. An optimal lossy approximation is chosen for each layer
9 by weighing the relative accuracy loss and FLOP reduction. On PASCAL VOC
10 2007 with the YOLO network, we show an end-to-end 2x speedup in a network
11 forward pass with a 5% drop in mAP that can be re-gained by finetuning, enabling
12 this network (and others like it) to be deployed in compute-constrained systems.

13 1 Introduction

14 At deploy time, the dollar cost of a production pipeline using a neural network is proportional to
15 the time required to execute a forward pass, which is proportional to the number of floating-point
16 operations (FLOPs) required. We want to run faster models to achieve higher throughput, lower
17 cost, better hardware utilization, and lower power, cooling, and CPU speed requirements. However,
18 training the smallest possible network to achieve the desired task is challenging.

19 Instead, we propose the DLA method to take a network that is slightly too slow and reduce the FLOP
20 count with minimal accuracy loss, reducing the need to train and re-train to find an appropriately-
21 sized network. We automatically select an appropriate singular value decomposition (SVD) to apply
22 to each weight tensor in the original network, replacing it with a set of weights that, when applied
23 sequentially, offer similar outputs but using fewer FLOPs. Intuitively, using low-rank decomposition
24 to determine which FLOPs to keep and which to approximate means that representation layers that
25 are redundant will have a lower-rank decomposition and yield more speedup for less accuracy loss
26 when approximated. The sensitivity of this approximation is controlled by a single parameter that
27 represents whether accuracy or speedup is more important. Because this method does not require
28 access to the original network training data, it can be used to speed up systems where networks are
29 a black-box or training data is confidential.

30 On a benchmark set of standard computer vision tasks and neural network architectures, we show
31 between a 1.5x and 2x speedup without additional re-training and minimal accuracy loss. In each
32 case, funetuning after applying DLA recovers lost accuracy.

33 This work extends that of Denton et al. [1], who use SVD and factorization approaches to approx-
34 imate weight layers, but rely on re-training after every layer has been approximated. DLA extends
35 this method by applying approximations in a zero-shot manner, without requiring re-training.

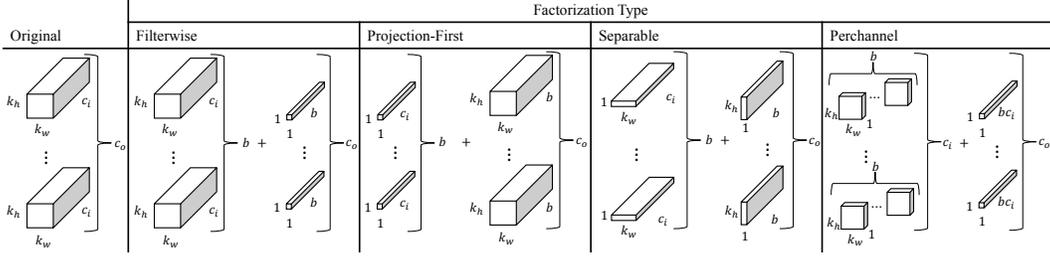


Figure 1: A visual representation of the 4 factorization methods used in DLA. c_i is the number of input channels, c_o is the number of output channels, k_w is the kernel width, k_h is the kernel height, and b is the factorization parameter.

2 Computational Model and Weight Tensors

According to the roofline computation model [11], network runtime is dominated by FLOP count and memory movement to and from the GPU. Weight layers fall in the FLOP-dominated regime because they have a high arithmetic intensity. DLA replaces FLOP-heavy operations with operations that have a lower FLOP count, pushing the layers' operations into the memory-limited regime.

The FLOP count is dominated by the duplication factor on the input blob size. A convolutional layer is a 4D tensor W in $\mathbb{R}^{c_o \times c_i \times k_h \times k_w}$, where c_o is the number of output feature maps, c_i is the number of input channels, and k_h and k_w are the kernel dimensions in y and x respectively. The convolution traverses an input blob in $\mathbb{R}^{h \times w}$ with stride s_h in y and s_w in x . In a grouped convolution, the input and output channels are broken into g groups and the 4D weight tensor is in $\mathbb{R}^{c_o \times \frac{c_i}{g} \times k_h \times k_w}$.

3 The Approximation Pipeline

Table 1: Properties of 4 types of SVD approximations, shown as a multiple of hw

	Filterwise	Projection-First	Separable	Perchannel
W_1 shape	$\mathbb{R}^{b \times c_i \times k_h \times k_w}$	$\mathbb{R}^{b \times c_i \times 1 \times 1}$	$\mathbb{R}^{b \times c_i \times k_h \times 1}$	$\mathbb{R}^{bc_o \times \frac{c_i}{c_o} \times k_h \times k_w}$
W_2 shape	$\mathbb{R}^{c_o \times b \times 1 \times 1}$	$\mathbb{R}^{c_o \times b \times k_h \times k_w}$	$\mathbb{R}^{c_o \times b \times 1 \times k_w}$	$\mathbb{R}^{c_o \times bc_o \times 1 \times 1}$
FLOPs for $W_1 + W_2$	$\frac{c_i b k_h k_w}{s_h s_w g} + \frac{bc_o}{s_h s_w}$	$c_i b + \frac{bc_o k_h k_w}{s_h s_w g}$	$\frac{c_i b k_h}{s_h s_w g} + \frac{bc_o k_w}{s_h s_w}$	$\frac{c_i b k_h k_w}{s_h s_w} + \frac{bc_o^2}{s_h s_w}$
A	$\frac{s_b^2}{\sum_{j=0}^{c_o} s_j^2}$	$\frac{s_b^2}{\sum_{j=0}^{c_o} s_j^2}$	$\frac{s_b^2}{\sum_{j=0}^{c_o} s_j^2}$	$\frac{1}{c_i} \sum_{c=0}^{c_i} \frac{s_c b^2}{\sum_{j=0}^{c_o} s_c j^2}$
R	$\frac{b}{c_o} + \frac{bg}{c_i k_h k_w}$	$\frac{bs_h s_w g}{c_o k_w k_h} + \frac{b}{c_i}$	$\frac{b}{c_o k_w} + \frac{bg}{c_i k_h}$	$\frac{bg}{c_o} + \frac{bgc_o}{c_i k_w k_h}$

An optimal approximation is chosen by calculating the runtime and accuracy loss from all possible decompositions (including chaining multiple approximations) and selecting the one with the highest score. The lossy approximations that are enumerated (both independently and chained together) are low-rank approximations of the original weight tensor W using SVD [4], resulting in two convolutions W_1 and W_2 that can be applied sequentially in place of W , shown visually in Figure 1. An **accuracy score** A (the percentage of variation explained by the approximation) and **runtime score** R (the FLOP reduction) is computed for each approximation. The operation with the highest combined score $pA + (1 - p)R$ is selected. A table of SVD decompositions and FLOP reductions for each type of approximation is shown in Table 1. When chaining approximations, R for is the ratio of the final output FLOPs to the FLOPs from W . A is the product of the accuracy scores for each approximation in the chain, since any error introduced by the first will be carried over to the next.

Table 2: DLA applied to 4 networks and datasets

Network / dataset	Runtime (ms)			Accuracy (top-1 or mAP)		
	Baseline	DLA	Speedup	Baseline	DLA	Finetuned
AlexNet [7] / CIFAR10 [6]	6.06	3.47	1.75x	70.3%	67.75%	68.8%
ResNet50 [3] / ImageNet2012 [9]	40.9	26.8	1.50x	72.3%	62.2%	68.6%
VGG16 [10] / ImageNet2012	78.6	45.2	1.75x	70.38%	59.60%	70.4%
YOLO [8] / VOC2007 [2]	63	31	2.00x	66.9	61.9	65.9

58 4 Experimental results

59 Table 2 shows between 1.5x and 2x runtime improvement, between 5 and 10% loss in accuracy, and
60 nearly full recovery of accuracy after finetuning on 4 standard networks.¹ FLOP reduction correlates
61 with the absolute speedup, with the exception of the ResNet50 network, as shown in Table 3. Table
62 4 shows that the input parameter p can be chosen based on the desired runtime / accuracy tradeoff.

Table 3: Speedup, FLOP reduction, and memory reduction

Network on dataset	Speedup	FLOP ↓	Memory ↓
AlexNet [7] on CIFAR10 [6]	1.75x	2.50x	1.10x
ResNet50 [3] on ImageNet2012 [9]	1.50x	1.20x	1.50x
VGG16 [10] on ImageNet2012	1.75x	1.70x	1.50x
YOLO [8] on VOC2007 [2]	2.00x	2.00x	1.60x

Table 4: Runtime and accuracy on YOLO

p	ms	mAP
baseline	63	66.9
0.9	57	66.9
0.8	54	66.0
0.7	45	65.4
0.6	41	65.0
0.5	32	61.9
0.4	27	50.0

63 Memory-limited layers such as fully-connected layers and 1×1 convolutions do not see as much
64 of an improvement in runtime from DLA (for example ResNet50 in Table 3). Additionally, pushing
65 beyond the 2x speedup observed on YOLO without significant accuracy loss is not possible with
66 DLA, because layers are moved from the FLOP-limited regime to the memory-limited regime so
67 more aggressive approximation trades off model accuracy for less runtime improvement.

68 5 Conclusion

69 After networks see diminishing returns with FLOP-reduction methods like DLA, it means that mem-
70 ory movement is the runtime bottleneck, and memory-reducing optimizations should be applied,
71 suggesting a future research direction into zero-shot memory reduction.

72 Deep Learning Approximation can be applied to an already-trained network to speed it up and incur
73 only a small amount of accuracy loss. Access to training data is not required and the techniques
74 are framework-agnostic, which means DLA can be used on black-box networks or in environments
75 where the training data cannot be accessed. The combination of approximations that best achieve
76 the desired accuracy loss is chosen for each layer through an exhaustive search. DLA can be com-
77 bined with other methods for speeding up neural networks. This runtime reduction can generate a
78 multiplier in cost reduction for production services that use neural networks. Any accuracy loss that
79 that was introduced can be recovered by fine-tuning or re-training the new resultant network.

¹Runtime was tested on a single input with Caffe [5] compiled with CUDA8 on a g2 cloud instance.

80 **References**

- 81 [1] Emily Denton, Wojciech Zaremba, Joan Bruna, Yann LeCun, and Rob Fergus. Exploiting linear structure
82 within convolutional networks for efficient evaluation. In *Proceedings of the 27th International Confer-*
83 *ence on Neural Information Processing Systems - Volume 1*, pages 1269–1277, 2014.
- 84 [2] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual ob-
85 ject classes (voc) challenge. *International Journal of Computer Vision*, 2007. URL [http://link.](http://link.springer.com/10.1007/s11263-009-0275-4)
86 [springer.com/10.1007/s11263-009-0275-4](http://link.springer.com/10.1007/s11263-009-0275-4).
- 87 [3] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
88 In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- 89 [4] M. Jaderberg, A. Vedaldi, and A. Zisserman. Speeding up convolutional neural networks with low rank
90 expansions. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2014.
- 91 [5] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio
92 Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. 2014.
- 93 [6] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research).
94 Technical report. URL <http://www.cs.toronto.edu/~kriz/cifar.html>.
- 95 [7] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional
96 neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in*
97 *Neural Information Processing Systems 25*, pages 1097–1105. 2012.
- 98 [8] Joseph Redmon and Ali Farhadi. Yolo9000: Better, faster, stronger. 2016.
- 99 [9] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang,
100 Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large
101 scale visual recognition challenge. *Int. J. Comput. Vision*, 10(1007):11263–015, December 2015.
- 102 [10] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition.
103 Technical report, 2014.
- 104 [11] Samuel Williams, Andrew Waterman, and David Patterson. Roofline: An insightful visual performance
105 model for multicore architectures. In *Commun. ACM*, pages 65–76, April 2009.