
Training Structured Efficient Convolutional Layers

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Typical recent neural network designs are primarily convolutional layers, but
2 the tricks enabling structured efficient linear layers (SELLs) have not yet been
3 adapted to the convolutional setting. We present a method to express the weight
4 tensor in a convolutional layer using diagonal matrices, discrete cosine transforms
5 (DCTs) and permutations that can be optimised using standard stochastic gradient
6 methods. A network composed of such structured efficient convolutional layers
7 (SECL) outperforms existing low-rank networks and demonstrates competitive
8 computational efficiency.

9 1 Introduction

10 Deep neural networks have evolved, and no longer contain the gigantic linear layers seen in Si-
11 monyan & Zisserman (2015), instead opting for large convolutional feature maps and increased
12 depth (Springenberg et al., 2014; He et al., 2016; Hu et al., 2018). This means that most of the
13 network’s parameters and multiply-add (Mult-Add) operations are used in these convolutions. We
14 present a method to reduce both, while keeping the network structure the same.

15 SELLs provide a framework to approximate linear layers that we adapt to make convolutions more
16 efficient. A convolution can be viewed as a matrix multiplication between the extracted patches from
17 the input tensor and the weights of the filters passed over the image. In this work, we show that this
18 matrix multiplication can be replaced with a structured efficient transformation based on the ACDC
19 layer (Moczulski et al., 2016). In practice, this gives us a specific parameterisation at training time,
20 which we describe in Section 3.

21 This structured efficient parameterisation of a convolutional layer consumes far fewer parameters
22 than a full convolutional layer, scaling as $O(N)$ versus $O(N^2)$ in the number of channels N . At
23 the same time, it can be implemented efficiently at test time using a fast DCT. In our experiments,
24 detailed in Section 4, we show that, using distillation (Crowley et al., 2018), a network composed of
25 primarily this type of convolution can learn to classify CIFAR-10 to 91.43% accuracy, while only
26 using 46,442 parameters. For comparison, this is a greater accuracy and fifty times fewer parameters
27 than the network presented in Hinton et al. (2016).

28 2 Background

29 SELLs are aimed at compressing the linear layers of convolutional networks. Typically, they are
30 composed as an operator Φ :

$$\mathbf{y} = \mathbf{x}\Phi(\mathbf{D}, \mathbf{P}, \mathbf{S}, \mathbf{B}) \quad (1)$$

31 In which, \mathbf{D} are diagonal matrices, \mathbf{P} are permutations, \mathbf{S} are sparse matrices and \mathbf{B} are bases, such
32 as Fourier, Cosine (Moczulski et al., 2016) or Hadamard (Yang et al., 2015; Ailon & Chazelle, 2009)

33 transforms. Using these component transformations, the resulting random projections can approxi-
 34 mate random matrices used in deep learning. In this work, we build on the ACDC SELL (Moczulski
 35 et al., 2016).

36 Our approach can be viewed as yielding a low-rank tensor to use in the convolutional layers. Previous
 37 efforts on low-rank convolutional networks have focused on transforming pre-trained networks (Jader-
 38 berg et al., 2014; Alvarez & Petersson, 2016; Denton et al., 2014; Lebedev et al., 2014) or training
 39 networks with appropriate regularisers (Alvarez & Salzmann, 2017; Wen et al., 2017).

40 Networks with low-rank constraints are markedly more difficult to train. In Garipov et al. (2016)
 41 the authors train such a network on CIFAR-10, but only achieve a $2\times$ compression rate over a
 42 convolutional network, and attain less than 90% accuracy. Other papers have focused on similar
 43 tensor decompositions; Su et al. (2018) obtain 91.28% accuracy compressing a ResNet-34. However,
 44 neither decomposition affects the number of Mult-Adds used at test time, whereas our method
 45 achieves a substantial reduction.

46 3 Structured Efficient Convolutional Layers

47 If we define a function to map the patches over which a convolution passes on an input tensor to
 48 rows of a matrix as \mathbf{F} , we can express convolution using a *kernel matrix* \mathbf{W} as $\mathbf{y} = \mathbf{F}^{-1}(\mathbf{F}(x)\mathbf{W})$.
 49 This is the common algorithm known as `im2col-gemm` (Chetlur et al., 2014). From this, we define a
 50 Structured Efficient Convolutional Layer (SECL) with the following parameterisation for \mathbf{W} :

$$\mathbf{W} = \prod_{l=1}^L \mathbf{A}_l \mathbf{C} \mathbf{D}_l \mathbf{C}^{-1} \mathbf{P} \quad (2)$$

51 Where \mathbf{A} and \mathbf{D} are diagonal matrices, \mathbf{C} and \mathbf{C}^{-1} are the forward and inverse DCTs and \mathbf{P} is a
 52 riffle shuffle.

53 Using the weight matrix in Equation 2 is equivalent to a stack of ACDC layers (Moczulski et al.,
 54 2016); the permutation being implemented by a riffle shuffle. A *riffle shuffle* is a fixed permutation,
 55 splitting the input in half and then interleaving the two halves; equivalent to a perfect riffle shuffle
 56 with a deck of cards (Gilbert, 1955). As described in Section 4, this was found to work as well as a
 57 fixed random permutation and can be evaluated much faster (Zhang et al., 2018).

58 Substituting this parameterisation into convolutional layers presents a problem: most kernel matrices
 59 are not square, with one exception. Kernel matrices in pointwise convolutions are square when the
 60 number of input channels matches the output. To increase the number of channels, we repeat the
 61 input along the channel dimension. As channels commonly increase in integer steps, this allows us to
 62 implement almost all pointwise convolutions.

63 Given a pointwise convolution, we can now implement a convolution with any kernel size by preceding
 64 the pointwise with a grouped convolution. This is known as a depthwise separable convolution and
 65 has been demonstrated as a substitute for convolution (Chollet, 2016). This can also be implemented
 66 using an ACDC parameterisation for each filter, but there is not much benefit, as shown in Section 4.

67 The motivation underlying ACDC layers comes from their complex equivalent; using Fourier trans-
 68 forms, \mathbf{F} , it is possible to show (Moczulski et al., 2016; Huhtanen, 2008) almost all matrices \mathbf{M} can
 69 be factored as:

$$\mathbf{M} = \left[\prod_{i=1}^{N-1} \mathbf{D}_{2i-1} \mathbf{R}_{2i} \right] \mathbf{D}_{2N-1} \quad (3)$$

70 Where \mathbf{D}_{2i-1} is a diagonal and \mathbf{R}_{2i} is composed of \mathbf{FDF}^{-1} .

71 However, machine learning systems typically operate using real numbers, leading to the decision to
 72 use the DCT. Despite the break in theory, the ACDC layer was found to work well as a replacement
 73 for the fully connected layers in CaffeNet (Moczulski et al., 2016).

Table 1: Results of training a Wide ResNet with SECL substituting the convolutional layers. WRN-SECL refers to a Wide ResNet using full rank grouped 3x3 convolutions, while WRN-SECL-LR refers to using ACDC-parameterised grouped 3x3 convolutions. Networks trained without distillation are reported under *Scratch*, while those trained using any form of distillation are under *Distilled*.

Model	Params	Mult-Adds	Scratch	Distilled
			Top 1	Top 1
WRN(40,2)	2.24M	328M	4.8%	–
WRN-SECL(40,2)	46.4K	33.2M	10.42%	8.57%
WRN-SECL-LR(40,2)	38.7K	33.2M	13.4%	10.85%
Belagiannis et al. (2018)	0.27M	–	–	8.08%
Hinton et al. (2016)	0.27M	–	–	8.88%
Su et al. (2018)	40K	–	–	8.72%

74 4 Experiments

75 We demonstrate the effectiveness of this compression strategy in Section 4.1 in experiments on
 76 CIFAR-10. PyTorch (Paszke et al.) was used to implement experiments. All code will be made
 77 available following the reviewing process.

78 **Linear Approximation** We compared the riffle shuffle to a fixed random permutation on the toy
 79 synthetic regression problem described in Section 6.1 of Moczulski et al. (2016). Both permutations
 80 converged to a final mean squared error of 0.02.

81 **Importance of Weight Decay** A small network based on the All Convolutional network (Springen-
 82 berg et al., 2014) was trained using Hyperband (Li et al., 2017) to tune the weight decay term, learning
 83 rate and minibatch size on the CIFAR-10 dataset (Krizhevsky, 2009). The networks converging well
 84 had weight decay around 10^{-5} , prompting us to use 8.8×10^{-6} in all experiments.

85 4.1 Convolutional Networks

86 To compare with contemporary works on compressed networks, we focus on a recent network
 87 architecture, but SECLs could be effectively substituted into any convolutional network. The network
 88 used in experiments was a Wide ResNet with depth 40 and width factor 2 (Zagoruyko & Komodakis,
 89 2016), and we train it on the CIFAR-10 dataset (Krizhevsky, 2009). We use 12 ACDC layers in each
 90 convolutional layer, to match the original paper (Moczulski et al., 2016).

91 The results are shown in Table 1. The results are most comparable to those of Su et al. (2018), in
 92 which the authors compress a network using a low-rank approximation. We performed distillation
 93 using attention-transfer (Zagoruyko & Komodakis, 2017) with the method described in Crowley et al.
 94 (2018). Each other paper reported as *distilled* in Table 1 implements its own form of distillation,
 95 detailed in each paper.

96 Each ACDC layer is expected to cost $4N + 5N \log_2(N)$ Mult-Adds at test time. However, we found
 97 that the early layers, where the dimension of the effective matrix multiplication is low, the Mult-Adds
 98 used by the unrolled ACDC layers is greater than just applying the parameterised filter tensor in a
 99 convolution. In this case, we assume that at test time the implementation of layers would depend on
 100 which would be cheaper.

101 5 Conclusion

102 It is well known that deep networks are overparameterised (Denil et al., 2013), but networks with
 103 low-rank weight matrices have failed to gain traction as they are typically difficult to train. Fortunately,
 104 for networks with structured efficient parameterisations we can use the same tools we would use
 105 to train a standard deep network, requiring a minimum of hyperparameter tuning. At test time,
 106 they can be implemented extremely efficiently using FFT-like algorithms or, as noted in Moczulski
 107 et al. (2016), by an optical processor (Saade et al., 2015). In this work, we have shown that these
 108 parameterisations aren’t restricted to linear layers, and can be applied to convolutional layers too;
 109 resulting in small, efficient neural networks.

References

- 110
- 111 Ailon, N. and Chazelle, B. The fast Johnson–Lindenstrauss transform and approximate nearest
112 neighbors. *SIAM Journal on computing*, 39(1):302–322, 2009.
- 113 Alvarez, J. M. and Salzmann, M. Compression-aware training of deep networks. In *Advances in*
114 *Neural Information Processing Systems*, 2017.
- 115 Alvarez, J. M. and Petersson, L. DecomposeMe: Simplifying convnets for end-to-end learning.
116 *arXiv:1606.05426*, 2016.
- 117 Belagiannis, V., Farshad, A., and Galasso, F. Adversarial network compression. *CoRR*,
118 abs/1803.10750, 2018. URL <http://arxiv.org/abs/1803.10750>.
- 119 Chetlur, S., Woolley, C., Vandermersch, P., Cohen, J., Tran, J., Catanzaro, B., and Shelhamer, E.
120 cuDNN: Efficient primitives for deep learning. *arXiv:1410.0759*, 2014.
- 121 Chollet, F. Xception: Deep learning with depthwise separable convolutions. *CoRR*, abs/1610.02357,
122 2016. URL <http://arxiv.org/abs/1610.02357>.
- 123 Crowley, E. J., Gray, G., and Storkey, A. Moonshine: Distilling with cheap convolutions. In *Advances*
124 *in Neural Information Processing Systems*, 2018.
- 125 Denil, M., Shakibi, B., Dinh, L., Ranzato, M., and de Freitas, N. Predicting parameters in deep
126 learning. In *Advances in Neural Information Processing Systems*, 2013.
- 127 Denton, E. L., Zaremba, W., Bruna, J., LeCun, Y., and Fergus, R. Exploiting linear structure within
128 convolutional networks for efficient evaluation. In *Advances in Neural Information Processing*
129 *Systems*, 2014.
- 130 Garipov, T., Podoprikin, D., Novikov, A., and Vetrov, D. P. Ultimate tensorization: compressing
131 convolutional and FC layers alike. *arXiv:1611.03214*, 2016.
- 132 Gilbert, E. Theory of shuffling. Technical report, Bell Labs, Murray Hill, New Jersey, U.S., 1955.
- 133 He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings*
134 *of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- 135 Hinton, G., Vinyals, O., and Dean, J. Distilling the knowledge in a neural network. *CoRR*,
136 abs/1503.02531, 2016. URL <http://arxiv.org/abs/1503.02531>.
- 137 Hu, J., Shen, L., and Sun, G. Squeeze-and-excitation networks. In *Proceedings of the IEEE*
138 *Conference on Computer Vision and Pattern Recognition*, 2018.
- 139 Huhtanen, M. Approximating ideal diffractive optical systems. *Journal of Mathematical Analysis*
140 *and Applications*, 345(1):53–62, 2008.
- 141 Jaderberg, M., Vedaldi, A., and Zisserman, A. Speeding up Convolutional Neural Networks with
142 Low Rank Expansions. In *British Machine Vision Conference*, 2014.
- 143 Krizhevsky, A. Learning multiple layers of features from tiny images. Master’s thesis, University of
144 Toronto, 2009.
- 145 Lebedev, V., Ganin, Y., Rakhuba, M., Oseledets, I., and Lempitsky, V. Speeding-up convolutional
146 neural networks using fine-tuned CP-decomposition. *arXiv:1412.6553*, 2014.
- 147 Li, L., Jamieson, K., DeSalvo, G., Rostamizadeh, A., and Talwalkar, A. Hyperband: A novel bandit-
148 based approach to hyperparameter optimization. *The Journal of Machine Learning Research*, 18
149 (1):6765–6816, 2017.
- 150 Moczulski, M., Denil, M., Appleyard, J., and de Freitas, N. ACDC: a structured efficient linear layer.
151 In *International Conference on Learning Representations*, 2016.
- 152 Paszke, A., Gross, S., Chintala, S., and Chanan, G. PyTorch: Tensors and dynamic neural networks
153 in Python with strong GPU acceleration. <https://github.com/pytorch/pytorch>. Accessed:
154 31st October 2017.
- 155 Saade, A., Caltagirone, F., Carron, I., Daudet, L., Drémeau, A., Gigan, S., and Krzakala, F. Random
156 projections through multiple optical scattering: Approximating kernels at the speed of light. *CoRR*,
157 abs/1510.06664, 2015. URL <http://arxiv.org/abs/1510.06664>.
- 158 Simonyan, K. and Zisserman, A. Very deep convolutional networks for large-scale image recognition.
159 In *International Conference on Learning Representations*, 2015.

- 160 Springenberg, J. T., Dosovitskiy, A., Brox, T., and Riedmiller, M. A. Striving for simplicity: The all
161 convolutional net. *arXiv:1412.6806*, 2014.
- 162 Su, J., Li, J., Bhattacharjee, B., and Huang, F. Tensorized Spectrum Preserving Compression for
163 Neural Networks. *arXiv:1805.10352*, 2018.
- 164 Wen, W., Xu, C., Wu, C., Wang, Y., Chen, Y., and Li, H. Coordinating filters for faster deep neural
165 networks. In *Proceedings of the IEEE International Conference on Computer Vision*, 2017.
- 166 Yang, Z., Moczulski, M., Denil, M., de Freitas, N., Smola, A., Song, L., and Wang, Z. Deep fried
167 convnets. In *Proceedings of the IEEE International Conference on Computer Vision*, 2015.
- 168 Zagoruyko, S. and Komodakis, N. Wide residual networks. In *British Machine Vision Conference*,
169 2016.
- 170 Zagoruyko, S. and Komodakis, N. Paying more attention to attention: Improving the performance
171 of convolutional neural networks via attention transfer. In *International Conference on Learning
172 Representations*, 2017.
- 173 Zhang, X., Zhou, X., Lin, M., and Sun, J. ShuffleNet: An extremely efficient convolutional neural
174 network for mobile devices. In *Proceedings of the IEEE Conference on Computer Vision and
175 Pattern Recognition*, 2018.