
Efficient Neural Network Compression via Transfer Learning for Industrial Optical Inspection

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this paper, we investigate learning the deep neural networks for automated optical
2 inspection in industrial manufacturing. Our preliminary result has shown the stun-
3 ning performance improvement by transfer learning from the completely dissimilar
4 source domain: ImageNet. Further study for demystifying this improvement shows
5 that the transfer learning produces a highly compressible network, which was not
6 the case for the network learned from scratch. The experimental result shows that
7 there is a negligible accuracy drop in the network learned by transfer learning until
8 it is compressed to 1/128 reduction of the number of convolution filters. This result
9 is contrary to the compression without transfer learning which loses more than 5%
10 accuracy at the same compression rate.

11 1 Introduction

12 Since every manufactured product must pass the inspection process of detecting defects on the product
13 surface, the fast and accurate inspection algorithm is essential in the manufacturing industry. Contrary
14 to the success in natural image detection using sufficient data [4, 7], the available images of product
15 surfaces, in particular those containing the patterns of defects, are insufficient for the reliable training
16 of the deep structured network. While training the network from scratch only produces 81.07%
17 performance, we previously reported the 99.90% performance when the whole network is fine-tuned
18 from a successfully learned network for imageNet classification [2]. The result is non-trivial because
19 the images in the ImageNet source domain do not overlap with the visual inspection images in the
20 target domain as shown in Figure 1.

21 In this paper, we investigate the compression of the learned network for fast inference. Once the
22 network is learned by transfer learning, the activation of the neurons is very sparse. Different from the
23 result in [2], the network trained from scratch can also achieve 99.78% accuracy with the extensively
24 augmented data and a long period of training. Surprisingly, however, although the performance of
25 both networks is similar, the network trained from scratch learns much denser features of the input
26 data than the network trained by transfer learning. We experimentally show that using standard
27 teacher-student model compression technique [3], the network trained by transfer learning can be
28 compressed to 1/128 reduction of the number of convolution filters with a negligible accuracy drop.
29 In contrast, the network trained from scratch loses more than 5% accuracy at the same compression
30 rate. Previous works on the network compression have focused on the methods of compression
31 [3, 5, 6, 9], but our work is the first report to the best of our knowledge that the *transfer learning* is
32 related to the network compression.

33 The rest of the paper is organized as follows. In Section 2, we explain the method of experiments and
34 show the compression result in Section 2.3. Finally, we conclude with a brief discussion in Section 3.

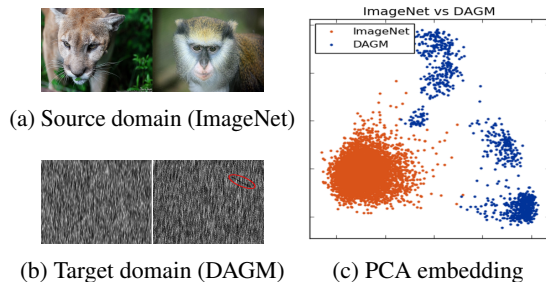


Figure 1: Sample images of the source and target domains, and their embedding result using PCA. Both domains are clearly dissimilar each other.

Table 1: Teacher network result

Method	Accuracy	Convergence time (SGD iterations)
Scratch (no aug)	81.07%	35,880
Scratch (aug)	99.78%	165,600
TL (no aug)	99.90%	2242

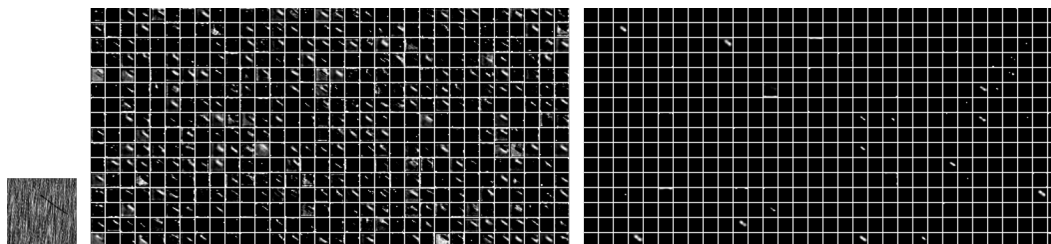


Figure 2: Output activations of the last convolutional layer of the Scratch (middle) and TL (right) teacher networks for a given input image (left).

35 2 Methods and experimental results

36 Following the neural network compression framework of [3], we first train the high performance
 37 deep teacher network and train a shallow student network to mimic it. To summarize the main results,
 38 the teacher network trained by transfer learning can be compressed to the shallow student network
 39 with 1/156 number of parameters and achieves 40× faster inference speed while showing less than
 40 1% accuracy drop (Table 2). Details of the experiments are explained below.

41 2.1 Target problem

42 The target problem is a texture inspection problem that the dataset is publicly accessible in the DAGM
 43 [1]. The dataset contains six patterns of texture with each pattern containing 1000 non-defective
 44 and 150 defective images, resulting in 6900 images of 12 classes (Figure 1b). The entire dataset is
 45 randomly divided into 80% for training (5520 images) and 20% for evaluation (1380 images).

46 2.2 Teacher networks: training from scratch vs transfer learning

47 We use the VGG16 [8] as a teacher network. We train the teacher network in two different ways with
 48 the same optimizer setting (SGD with batch size 32 and learning rate 10^{-3}). In the first case, we
 49 randomly initialize all weights in the network and train it directly on the target problem from scratch.
 50 We denote this deep teacher network as a *Scratch*. Since the number of training data is insufficient,
 51 the accuracy of the network no longer increases at 81.07% without any data augmentation. On the
 52 other hand, if we augment the training data by adding rotated and flipped images, the network can be
 53 successfully trained to achieve 99.78% accuracy after 165k iterations of training (Table 1).

54 For the second case of the teacher network, we apply transfer learning to train the network. We
 55 denote this teacher network as a *TL*. The source domain of the transfer learning is the ImageNet data
 56 [7] (Figure 1a) where the images bear almost no similarity to the target inspection data (Figure 1b).
 57 We can clearly see the dissimilarity between each domain by embedding features extracted from
 58 both datasets using principal component analysis (PCA) (Figure 1c). We follow the transfer learning
 59 method of [2] that weights of the target network initialized with the weights of the source network
 60 only up to convolutional layers. Then, the network is fine-tuned using target data without any data
 61 augmentation. Even though the source and target domains are completely dissimilar from each other,
 62 the network trained by transfer learning converge to 99.90% accuracy 70× faster than training from
 63 scratch (Table 1). In addition, when we visualize the activation of the last convolutional layer for

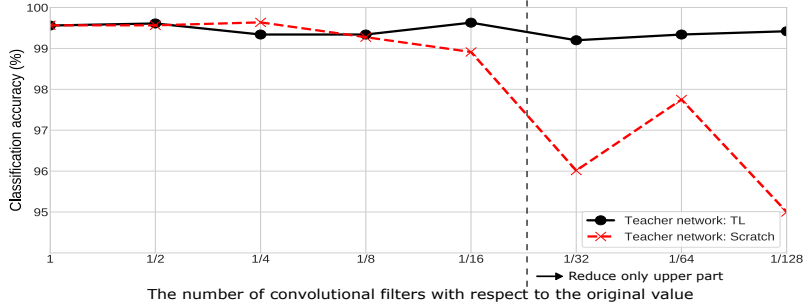


Figure 3: Neural network compression results (Scratch teacher network vs TL teacher network).

Table 2: Summary of neural network compression results

Model	Number of parameters	Accuracy (%)	Inference speed (ms/image)
Teacher Network (TL)	134 million	99.90	57.43
Student Network (TL)	858,380	99.42	1.46
Student Network (Scratch)	858,380	94.92	1.46

64 a given input image, we can see that the TL teacher network learns much sparser features of the
 65 input data than the Scratch teacher network (Figure 2). The sparsity of activation in the TL teacher
 66 network suddenly increases sharply from the eighth convolutional layer and reaches 99.10% at the
 67 last convolutional layer.

68 2.3 Neural network compression: *Scratch* teacher network vs *TL* teacher network

69 The student network is trained to learn the input-output mapping function of the teacher network
 70 by regressing the logits output of the teacher network [3]. For training the student network, we
 71 use the data same as the teacher network training set with the augmented version. Focusing on
 72 accelerating the inference time, we design the shallow student network as following rules. We remove
 73 the first fully-connected layer to reduce the number of parameters in the network. In addition, we
 74 eliminate eight of the 13 convolutional layers that account for a large part of the computational
 75 cost. The number of filters in remaining convolutional layers is denoted as $N_1, N_2, N_3, N_4,$ and N_5
 76 respectively. Then, we gradually reduce the number of filters in all convolutional layers by half from
 77 the original value ($N_1, N_2, N_3, N_4, N_5) = (64, 128, 256, 512, 512)$. The accuracy of both teacher
 78 networks is maintained until we reduce the number of all convolutional filters to 1/16. However, when
 79 it is reduced to 1/32, a sudden accuracy drop occurs to 82% accuracy in both networks. Using the fact
 80 that the sparsity of the activation grows rapidly from the eighth convolutional layer of the TL teacher
 81 network, we further reduce only upper part of the student network (N_4, N_5). The TL teacher network
 82 has a negligible accuracy drop until the upper part is compressed to 1/128. In contrast, compressing
 83 the Scratch teacher network loses more than 5% accuracy at the same compression point (Figure 3).
 84 Combining the results, the final compressed student model is $(N_1, N_2, N_3, N_4, N_5) = (4, 8, 16, 4, 4)$.
 85 The student network compressed from TL teacher network has 1/156 number of parameters and
 86 achieves 40× faster inference speed while only 0.48% accuracy drop (Table 2).

87 3 Conclusion

88 Through experiments using inspection data, we show that the high performance, lightweight shallow
 89 network can be obtained for optical inspection via compressing the deep network trained by transfer
 90 learning. When the performance of two networks is similar, sparse features learned in the network
 91 seems to play a key role in model compression. This result proposes that in addition to the compression
 92 method, what the deep network has learned is also important to the model compression framework.

93 **References**

- 94 [1] Weakly supervised learning for industrial optical inspection. [https://hci.iwr.](https://hci.iwr.uni-heidelberg.de/node/3616)
95 [uni-heidelberg.de/node/3616](https://hci.iwr.uni-heidelberg.de/node/3616). Accessed: 2018-10-20.
- 96 [2] ———. Transfer learning for automated optical inspection. In *International Joint Conference*
97 *on Neural Networks (IJCNN), 2017*, pages 2517–2524, 2017.
- 98 [3] Jimmy Ba and Rich Caruana. Do deep nets really need to be deep? In *Advances in neural*
99 *information processing systems*, pages 2654–2662, 2014.
- 100 [4] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman.
101 The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):
102 303–338, 2010.
- 103 [5] Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural net-
104 works with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*,
105 2015.
- 106 [6] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network.
107 *arXiv preprint arXiv:1503.02531*, 2015.
- 108 [7] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
109 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei.
110 ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*
111 (*IJCV*), 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- 112 [8] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale
113 image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 114 [9] Xiyu Yu, Tongliang Liu, Xinchao Wang, and Dacheng Tao. On compressing deep models by low
115 rank and sparse decomposition. In *Proceedings of the IEEE Conference on Computer Vision and*
116 *Pattern Recognition*, pages 7370–7379, 2017.