

# Similarity-Preserving Knowledge Distillation

Frederick Tung<sup>1,2</sup> and Greg Mori<sup>1,2</sup>  
<sup>1</sup>Simon Fraser University <sup>2</sup>Borealis AI  
 ftung@sfu.ca, mori@cs.sfu.ca

## Abstract

Knowledge distillation is a widely applicable technique for training a student neural network under the guidance of a trained teacher network. For example, in neural network compression, a high-capacity teacher is distilled to train a compact student; in privileged learning, a teacher trained with privileged data is distilled to train a student without access to that data. The distillation loss determines how a teacher’s knowledge is captured and transferred to the student. In this paper, we propose a new form of knowledge distillation loss that is inspired by the observation that semantically similar inputs tend to elicit similar activation patterns in a trained network. Similarity-preserving knowledge distillation guides the training of a student network such that input pairs that produce similar (dissimilar) activations in the teacher network produce similar (dissimilar) activations in the student network. In contrast to previous distillation methods, the student is not required to mimic the representation space of the teacher, but rather to preserve the pairwise similarities in its own representation space. Experiments on three public datasets demonstrate the potential of our approach.

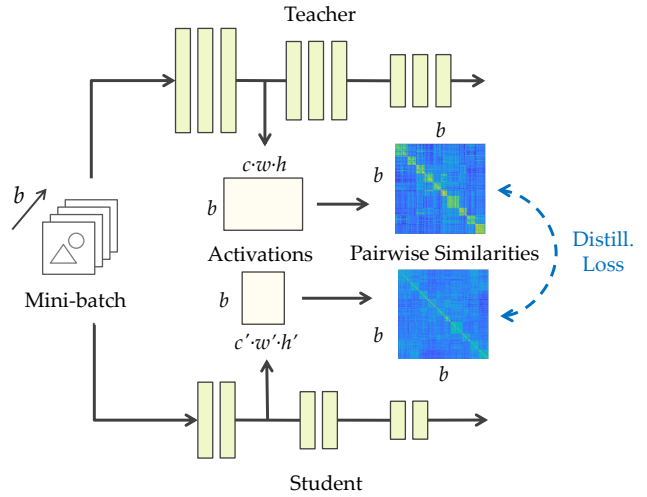


Figure 1. Similarity-preserving knowledge distillation guides the training of a student network such that input pairs that produce similar (dissimilar) activations in the pre-trained teacher network produce similar (dissimilar) activations in the student network. Given an input mini-batch of  $b$  images, we derive  $b \times b$  pairwise similarity matrices from the activation maps, and compute a distillation loss on the matrices produced by the student and the teacher.

## 1. Introduction

Deep neural networks are being used to solve an increasingly wide array of computer vision problems. While the general trend in deep learning is towards deeper, wider, and more complex networks, deploying deep learning solutions in the real world requires us to consider the computational cost. A mobile robot or self-driving vehicle, for example, has limited memory and power. Even when resources are abundant, such as when a vision system is hosted in the cloud, more resource-efficient deep networks mean more clients can be served at a lower cost. When performing transfer learning in the real world, data privilege and privacy issues may restrict access to data in the source domain. It may be necessary to transfer the knowledge of a network trained on the source domain assuming access only to training data from the target task domain.

Knowledge distillation is a general technique for supervising the training of “student” neural networks by capturing and transferring the knowledge of trained “teacher” networks. While originally motivated by the task of neural network compression for resource-efficient deep learning [12], knowledge distillation has found wider applications in such areas as privileged learning [21], adversarial defense [25], and learning with noisy data [19]. Knowledge distillation is conceptually simple: it guides the training of a student network with an additional distillation loss that encourages the student to mimic some aspect of a teacher network. Intuitively, the trained teacher network provides a richer supervisory signal than the data supervision (e.g. annotated class labels) alone.

The conceptual simplicity of knowledge distillation belies the fact that *how* to best capture the knowledge of the teacher to train the student (i.e. how to define the distillation

loss) remains an open question. In traditional knowledge distillation [12], the softened class scores of the teacher are used as the extra supervisory signal: the distillation loss encourages the student to mimic the scores of the teacher. FitNets [31] extend this idea by adding hints to guide the training of intermediate layers. In flow-based knowledge distillation [38], the extra supervisory signal comes from the inter-layer “flow” – how features are transformed between layers. The distillation loss encourages the student to mimic the teacher’s flow matrices, which are derived from the inner product between feature maps in two layers, such as the first and last layers in a residual block. In attention transfer [41], the supervisory signal for knowledge distillation is in the form of spatial attention. Spatial attention maps are computed by summing the squared activations along the channel dimension. The distillation loss encourages the student to produce similar normalized spatial attention maps as the teacher, intuitively paying attention to similar parts of the image as the teacher.

In this paper, we present a novel form of knowledge distillation that is inspired by the observation that semantically similar inputs tend to elicit similar activation patterns in a trained neural network. Similarity-preserving knowledge distillation guides the training of a student network such that input pairs that produce similar (dissimilar) activations in the trained teacher network produce similar (dissimilar) activations in the student network. Figure 1 shows the overall procedure. Given an input mini-batch of  $b$  images, we compute pairwise similarity matrices from the output activation maps. The  $b \times b$  matrices encode the similarities in the activations of the network as elicited by the images in the mini-batch. Our distillation loss is defined on the pairwise similarity matrices produced by the student and the teacher.

To support the intuition of our distillation loss, Figure 2 visualizes the average activation of each channel in the last convolutional layer of a WideResNet-16-2 teacher network (we adopt the standard notation WideResNet- $d$ - $k$  to refer to a wide residual network [40] with depth  $d$  and width multiplier  $k$ ), on the CIFAR-10 test images. We can see that images from the same object category tend to activate similar channels in the trained network. The similarities in activations across different images capture useful semantics learned by the teacher network. We study whether these similarities provide an informative supervisory signal for knowledge distillation.

The contributions of this paper are:

- We introduce *similarity-preserving knowledge distillation*, a novel form of knowledge distillation that uses the pairwise activation similarities within each input mini-batch to supervise the training of a student network with a trained teacher network.

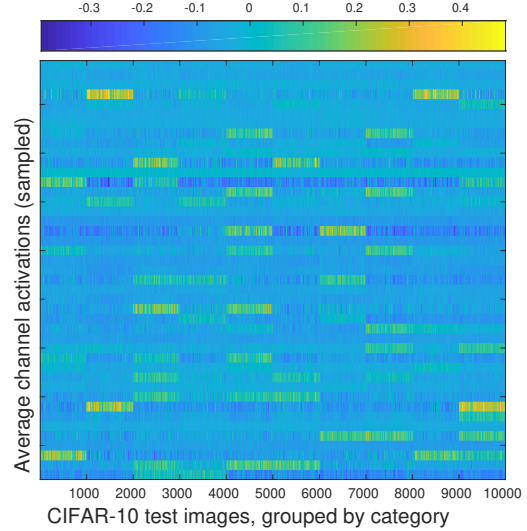


Figure 2. Semantically similar inputs tend to elicit similar activation patterns in a trained neural network. This visualization shows channel-wise average activations sampled from the last convolutional layer of a WideResNet-16-2 network on the CIFAR-10 test images. Activation patterns are largely consistent within the same category (e.g. columns 1 to 1000) and distinctive across different categories (e.g. columns 1 to 1000 vs. columns 1001 to 2000).

- We experimentally validate our approach on three public datasets. Our experiments show the potential of similarity-preserving knowledge distillation, not only for improving the training outcomes of student networks, but also for complementing traditional methods for knowledge distillation.

## 2. Method

The goal of knowledge distillation is to train a student network under the guidance of a trained teacher network, which acts as an extra source of supervision. For example, in neural network compression, the student network is computationally cheaper than the teacher: it may be shallower, thinner, or composed of cheaper operations. The trained teacher network provides additional semantic knowledge beyond the usual data supervision (e.g. the usual one-hot vectors for classification). The challenge is to determine how to encode and transfer the teacher’s knowledge such that student performance is maximized.

In traditional knowledge distillation [12], knowledge is encoded and transferred in the form of softened class scores. The total loss for training the student is given by

$$\mathcal{L} = (1 - \alpha)\mathcal{L}_{\text{CE}}(\mathbf{y}, \sigma(\mathbf{z}_S)) + 2\alpha T^2 \mathcal{L}_{\text{CE}}(\sigma(\frac{\mathbf{z}_S}{T}), \sigma(\frac{\mathbf{z}_T}{T})), \quad (1)$$

where  $\mathcal{L}_{\text{CE}}(\cdot, \cdot)$  denotes the cross-entropy loss,  $\sigma(\cdot)$  denotes the softmax function,  $\mathbf{y}$  is the one-hot vector indicating the ground truth class,  $\mathbf{z}_S$  and  $\mathbf{z}_T$  are the output logits of the

WideResNet-16-1  
(0.2M params)

WideResNet-40-2  
(2.2M params)

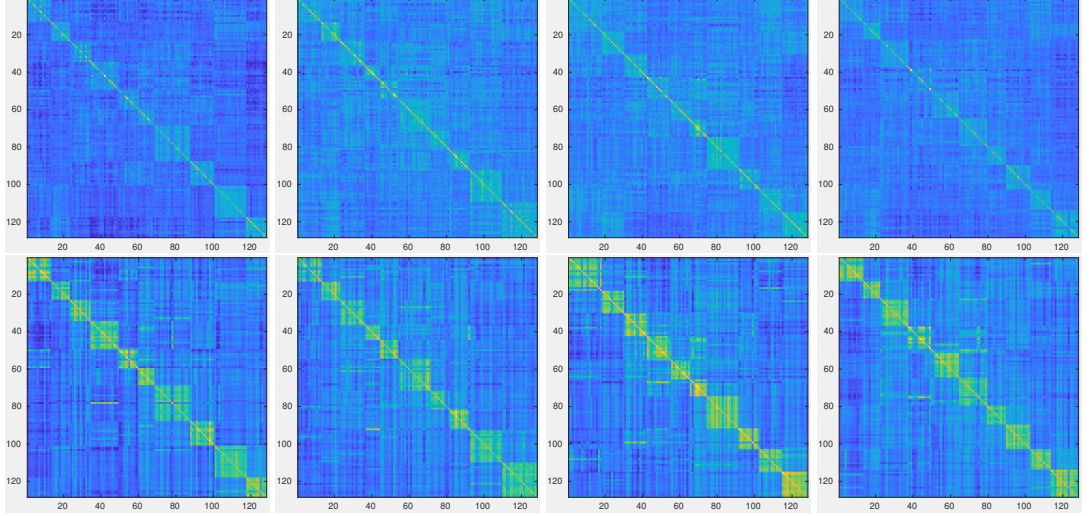


Figure 3. Activation similarity matrices  $G$  (Eq. 2) produced by trained WideResNet-16-1 and WideResNet-40-2 networks on sample CIFAR-10 test batches. Each column shows a single batch with inputs grouped by ground truth class along each axis (batch size = 128). Brighter colors indicate higher similarity values. The blockwise patterns indicate that the elicited activations are mostly similar for inputs of the same class, and different for inputs across different classes. Our distillation loss (Eq. 4) encourages the student network to produce  $G$  matrices closer to those produced by the teacher network.

student and teacher networks, respectively,  $T$  is a temperature hyperparameter, and  $\alpha$  is a balancing hyperparameter. The first term in Eq. 1 is the usual cross-entropy loss defined using data supervision (ground truth labels), while the second term encourages the student to mimic the softened class scores of the teacher.

Recall from the introduction and Figure 2 that semantically similar inputs tend to elicit similar activation patterns in a trained neural network. In Figure 2, we can observe that activation patterns are largely consistent within the same object category and distinctive across different categories. Might the correlations in activations encode useful teacher knowledge that can be transferred to the student? Our hypothesis is that, if two inputs produce highly similar activations in the teacher network, it is beneficial to guide the student network towards a configuration that also results in the two inputs producing highly similar activations in the student. Conversely, if two inputs produce dissimilar activations in the teacher, we want these inputs to produce dissimilar activations in the student as well.

Given an input mini-batch, denote the activation map produced by the teacher network  $T$  at a particular layer  $l$  by  $A_T^{(l)} \in \mathbf{R}^{b \times c \times h \times w}$ , where  $b$  is the batch size,  $c$  is the number of output channels, and  $h$  and  $w$  are spatial dimensions. Let the activation map produced by the student network  $S$  at a corresponding layer  $l'$  be given by  $A_S^{(l')} \in \mathbf{R}^{b \times c' \times h' \times w'}$ . Note that  $c$  does not necessarily have to equal  $c'$ , and likewise for the spatial dimensions. Similar to attention transfer [41], the corresponding layer  $l'$  can be the layer at the same depth as  $l$  if the student and teacher share the same

depth, or the layer at the end of the same block if the student and teacher have different depths. To guide the student towards the activation correlations induced in the teacher, we define a distillation loss that penalizes differences in the L2-normalized outer products of  $A_T^{(l)}$  and  $A_S^{(l')}$ . First, let

$$G_T^{(l)} = \frac{Q_T^{(l)} \cdot Q_T^{(l)\top}}{\|Q_T^{(l)} \cdot Q_T^{(l)\top}\|_2}, \quad (2)$$

where  $Q_T^{(l)} \in \mathbf{R}^{b \times chw}$  is a reshaping of  $A_T^{(l)}$ , and therefore  $G_T^{(l)}$  is a  $b \times b$  matrix. Intuitively, entry  $(i, j)$  in  $G_T^{(l)}$  encodes the similarity of the activations at this teacher layer elicited by the  $i$ th and  $j$ th images in the mini-batch. Analogously, let

$$G_S^{(l')} = \frac{Q_S^{(l')} \cdot Q_S^{(l')\top}}{\|Q_S^{(l')} \cdot Q_S^{(l')\top}\|_2}, \quad (3)$$

where  $Q_S^{(l')} \in \mathbf{R}^{b \times c' h' w'}$  is a reshaping of  $A_S^{(l')}$ , and  $G_S^{(l')}$  is a  $b \times b$  matrix. We define the similarity-preserving knowledge distillation loss as:

$$\mathcal{L}_{\text{SP}}(G_T, G_S) = \frac{1}{b^2} \sum_{(l, l') \in \mathcal{I}} \|G_T^{(l)} - G_S^{(l')}\|_F^2, \quad (4)$$

where  $\mathcal{I}$  collects the  $(l, l')$  layer pairs (e.g. layers at the end of the same block, as discussed above) and  $\|\cdot\|_F$  is the Frobenius norm. Eq. 4 is a summation, over all  $(l, l')$  pairs, of the mean element-wise squared difference between the  $G_T^{(l)}$  and  $G_S^{(l')}$  matrices. Finally, we define the total loss for training the student network as:

$$\mathcal{L} = \mathcal{L}_{\text{CE}}(\mathbf{y}, \sigma(\mathbf{z}_S)) + \gamma \mathcal{L}_{\text{SP}}(G_T, G_S), \quad (5)$$

Group name	Output size	WideResNet-16- $k$	WideResNet-40- $k$
conv1	$32 \times 32$	$3 \times 3, 16$	$3 \times 3, 16$
conv2	$32 \times 32$	$\begin{bmatrix} 3 \times 3, 16k \\ 3 \times 3, 16k \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 16k \\ 3 \times 3, 16k \end{bmatrix} \times 6$
conv3	$16 \times 16$	$\begin{bmatrix} 3 \times 3, 32k \\ 3 \times 3, 32k \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 32k \\ 3 \times 3, 32k \end{bmatrix} \times 6$
conv4	$8 \times 8$	$\begin{bmatrix} 3 \times 3, 64k \\ 3 \times 3, 64k \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64k \\ 3 \times 3, 64k \end{bmatrix} \times 6$
	$1 \times 1$	average pool, 10-d fc, softmax	

Table 1. Structure of WideResNet networks used in CIFAR-10 experiments. Downsampling is performed by strided convolutions in the first layers of conv3 and conv4.

where  $\gamma$  is a balancing hyperparameter.

Figure 3 visualizes the  $G$  matrices for several batches in the CIFAR-10 test set. The top row is produced by a trained WideResNet-16-1 network, consisting of 0.2M parameters, while the bottom row is produced by a trained WideResNet-40-2 network, consisting of 2.2M parameters. In both cases, activations are collected from the last convolution layer. Each column represents a single batch, which is identical for both networks. The images in each batch have been grouped by their ground truth class for easier interpretability. The  $G$  matrices in both rows show a distinctive block-wise pattern, indicating that the activations at the last layer of these networks are largely similar within the same class and dissimilar across different classes (the blocks are differently sized because each batch has an unequal number of test samples from each class). Moreover, the blockwise pattern is more distinctive for the WideResNet-40-2 network, reflecting the higher capacity of this network to capture the semantics of the dataset. Intuitively, Eq. 4 pushes the student network towards producing  $G$  matrices closer to those produced by the teacher network.

**Differences from previous approaches.** The similarity-preserving knowledge distillation loss (Eq. 4) is defined in terms of activations instead of class scores as in traditional distillation [12]. Activations are also used to define the distillation losses in FitNets [31], flow-based distillation [38], and attention transfer [41]. However, a key difference is that these previous distillation methods encourage the student to mimic different aspects of the representation space of the teacher. Our method is a departure from this common approach in that it aims to preserve the pairwise activation similarities of input samples. Its behavior is unchanged by a rotation of the teacher’s representation space, for example. In similarity-preserving knowledge distillation, the student is not required to be able to express the representation space of the teacher, as long as pairwise similarities in the teacher space are well preserved in the student space.

### 3. Experiments

We now turn to the experimental validation of our distillation approach on three public datasets. We start with CIFAR-10 as it is a commonly adopted dataset for comparing distillation methods, and its relatively small size allows multiple student and teacher combinations to be evaluated. We then consider the task of transfer learning, and show how distillation and fine-tuning can be combined to perform transfer learning on a texture dataset with limited training data. Finally, we report results on the larger CINIC-10 dataset.

#### 3.1. CIFAR-10

CIFAR-10 consists of 50,000 training images and 10,000 testing images at a resolution of  $32 \times 32$ . The dataset covers ten object classes, with each class having an equal number of images. We conducted experiments using wide residual networks (WideResNets) [40] following [4, 41]. Table 1 summarizes the structure of the networks. We adopted the standard protocol [40] for training wide residual networks on CIFAR-10 (SGD with Nesterov momentum; 200 epochs; batch size of 128; and an initial learning rate of 0.1, decayed by a factor of 0.2 at epochs 60, 120, and 160). We applied the standard horizontal flip and random crop data augmentation. We performed baseline comparisons with respect to traditional knowledge distillation (softened class scores) and attention transfer. For traditional knowledge distillation [12], we set  $\alpha = 0.9$  and  $T = 4$  following the CIFAR-10 experiments in [4, 41]. Attention transfer losses were applied for each of the three residual block groups. We set the weight of the distillation loss in attention transfer and similarity-preserving distillation by held-out validation on a subset of the training set ( $\beta = 1000$  for attention transfer,  $\gamma = 3000$  for similarity-preserving distillation).

Table 2 shows our results experimenting with several student-teacher network pairs. We tested cases in which the student and teacher networks have the same width but different depth (WideResNet-16-1 student with WideResNet-40-1 teacher; WideResNet-16-2 student with WideResNet-40-2 teacher), the student and teacher networks have the same depth but different width (WideResNet-16-1 student with WideResNet-16-2 teacher; WideResNet-16-2 student with WideResNet-16-8 teacher), and the student and teacher have different depth and width (WideResNet-40-2 student with WideResNet-16-8 teacher). In all cases, transferring the knowledge of the teacher network using similarity-preserving distillation improved student training outcomes. Compared to conventional training with data supervision (i.e. one-hot vectors), the student network consistently obtained lower median error, from 0.5 to 1.2 absolute percentage points, or 7% to 14% relative, with no additional network parameters or operations. Similarity-preserving distillation also performed favorably with respect to the tra-

Student	Teacher	Student	KD [12]	AT [41]	SP (ours)	Teacher
WideResNet-16-1 (0.2M)	WideResNet-40-1 (0.6M)	8.74	8.48	8.30	<b>8.13</b>	6.51
WideResNet-16-1 (0.2M)	WideResNet-16-2 (0.7M)	8.74	7.94	8.28	<b>7.52</b>	6.07
WideResNet-16-2 (0.7M)	WideResNet-40-2 (2.2M)	6.07	6.00	5.89	<b>5.52</b>	5.18
WideResNet-16-2 (0.7M)	WideResNet-16-8 (11.0M)	6.07	5.62	5.47	<b>5.34</b>	4.24
WideResNet-40-2 (2.2M)	WideResNet-16-8 (11.0M)	5.18	4.86	<b>4.47</b>	4.55	4.24

Table 2. Experiments on CIFAR-10 with three different knowledge distillation losses: softened class scores (traditional KD), attention transfer (AT), and similarity preserving (SP). The median error over five runs is reported, following the protocol in [40, 41]. The best result for each experiment is shown in bold. Brackets indicate model size in number of parameters.

Output size	MobileNet- $k$
$112 \times 112$	$3 \times 3, 32k$
$112 \times 112$	$3 \times 3$ dw, $32k$ $1 \times 1, 64k$
$56 \times 56$	$3 \times 3$ dw, $64k$ $1 \times 1, 128k$
$56 \times 56$	$3 \times 3$ dw, $128k$ $1 \times 1, 128k$
$28 \times 28$	$3 \times 3$ dw, $128k$ $1 \times 1, 256k$
$28 \times 28$	$3 \times 3$ dw, $256k$ $1 \times 1, 256k$
$14 \times 14$	$3 \times 3$ dw, $256k$ $1 \times 1, 512k$
$14 \times 14$	$3 \times 3$ dw, $512k$ $1 \times 1, 512k$ } $\times 5$
$7 \times 7$	$3 \times 3$ dw, $512k$ $1 \times 1, 1024k$
$7 \times 7$	$3 \times 3$ dw, $1024k$ $1 \times 1, 1024k$
$1 \times 1$	average pool, 47-d fc, softmax

Table 3. Structure of MobileNet networks used in transfer learning experiments. ‘dw’ denotes depthwise convolution. Downsampling is performed by strided  $3 \times 3$  depthwise convolutions.

ditional (softened class scores) and attention transfer baselines, achieving the lowest error in four of the five cases. This validates our intuition that the activation similarities across images encode useful semantics learned by the teacher network, and provide an effective supervisory signal for knowledge distillation.

While we have presented these results from the perspective of improving the training of a student network, it is also possible to view the results from the perspective of the teacher network. Our results suggest the potential for using similarity-preserving distillation to compress large networks into more resource-efficient ones with minimal accuracy loss. In the fifth test, for example, the knowledge of a trained WideResNet-16-8 network, which contains 11.0M parameters, is distilled into a much smaller WideResNet-

Output size	MobileNetV2- $k$
$112 \times 112$	$3 \times 3, 32k$
$112 \times 112$	bottleneck( $t = 1, c = 16k, n = 1$ )
$56 \times 56$	bottleneck( $t = 6, c = 24k, n = 2$ )
$28 \times 28$	bottleneck( $t = 6, c = 32k, n = 3$ )
$14 \times 14$	bottleneck( $t = 6, c = 64k, n = 4$ )
$14 \times 14$	bottleneck( $t = 6, c = 96k, n = 3$ )
$7 \times 7$	bottleneck( $t = 6, c = 160k, n = 3$ )
$7 \times 7$	bottleneck( $t = 6, c = 320k, n = 1$ )
$7 \times 7$	$1 \times 1, 1280k$
$1 \times 1$	average pool, 47-d fc, softmax

Table 4. Structure of MobileNetV2 networks used in transfer learning experiments. The notation ‘bottleneck( $t, c, n$ )’ denotes a group of bottleneck residual blocks with expansion factor  $t$ ,  $c$  output channels, and  $n$  repeated blocks. Downsampling is performed by strided  $3 \times 3$  depthwise convolution in the first block of a group.

40-2 network, which contains only 2.2M parameters. This is a  $5\times$  compression rate with only 0.3% loss in accuracy, using off-the-shelf PyTorch without any specialized hardware or software.

The above similarity-preserving distillation results were produced using only the activations collected from the last convolution layers of the student and teacher networks. We also experimented with using the activations at the end of each WideResNet block, but found no improvement in performance. We therefore used only the activations at the final convolution layers in the subsequent experiments. Activation similarities may be less informative in the earlier layers of the network because these layers encode more generic features, which tend to be present across many images. Progressing deeper in the network, the channels encode increasingly specialized features, and the activation patterns of semantically similar images become more distinctive.

### 3.2. Transfer learning combining distillation with fine-tuning

In this section, we explore a common transfer learning scenario in computer vision. Suppose we are faced with a novel recognition task in a specialized image domain with



Student	Teacher	Student	AT [41]	SP (win:loss)	Teacher
MobileNet-0.25 (0.2M)	MobileNet-0.5 (0.8M)	42.45	42.39	<b>41.30</b> (7:3)	36.76
MobileNet-0.25 (0.2M)	MobileNet-1.0 (3.3M)	42.45	41.89	<b>41.76</b> (5:5)	34.10
MobileNet-0.5 (0.8M)	MobileNet-1.0 (3.3M)	36.76	35.61	<b>35.45</b> (7:3)	34.10
MobileNetV2-0.35 (0.5M)	MobileNetV2-1.0 (2.2M)	41.25	41.60	<b>40.29</b> (8:2)	36.62
MobileNetV2-0.35 (0.5M)	MobileNetV2-1.4 (4.4M)	41.25	41.04	<b>40.43</b> (8:2)	35.35
MobileNetV2-1.0 (2.2M)	MobileNetV2-1.4 (4.4M)	36.62	36.33	<b>35.61</b> (8:2)	35.35

Table 5. Transfer learning experiments on the describable textures dataset with attention transfer (AT) and similarity preserving (SP) knowledge distillation. The median error over the ten standard splits is reported. The best result for each experiment is shown in bold. The (win:loss) notation indicates the number of splits in which SP outperformed AT. The (\*M) notation indicates model size in number of parameters.

limited training data. A natural strategy to adopt is to transfer the knowledge of a network pre-trained on ImageNet (or another suitable large-scale dataset) to the new recognition task by fine-tuning. Here, we combine knowledge distillation with fine-tuning: we initialize the student network with source domain (in this case, ImageNet) pretrained weights, and then fine-tune the student to the target domain using both distillation and cross-entropy losses (Eq. 5).

We analyzed this scenario using the describable textures dataset [3], which is composed of 5,640 images covering 47 texture categories. Image sizes range from 300x300 to 640x640. We applied ImageNet-style data augmentation with horizontal flipping and random resized cropping during training. At test time, images were resized to 256x256 and center cropped to 224x224 for input to the networks. For evaluation, we adopted the standard ten training-validation-testing splits. To demonstrate the versatility of our method on different network architectures, and in particular its compatibility with mobile-friendly architectures, we experimented with variants of MobileNet [13] and MobileNetV2 [32]. Tables 3 and 4 summarize the structure of the networks.

We compared with an attention transfer baseline. Softened class score based distillation is not directly comparable in this setting because the classes in the source and target domains are disjoint. The teacher would first have to be fine-tuned to the target domain, which significantly increases training time and may not be practical when employing expensive teachers or transferring to large datasets. Similarity-preserving distillation can be applied directly to train the student, without first fine-tuning the teacher, since it aims to preserve similarities instead of mimicking the teacher’s representation space. We set the hyperparameters for attention transfer and similarity-preserving distillation by held-out validation on the ten standard splits. All networks were trained using SGD with Nesterov momentum, a batch size of 96, and for 60 epochs with an initial learning rate of 0.01 reduced to 0.001 after 30 epochs.

Table 5 shows that similarity-preserving distillation can effectively transfer knowledge across different domains.

For all MobileNet and MobileNetV2 student-teacher pairs tested, applying similarity-preserving distillation during fine-tuning resulted in lower median student error than fine-tuning without distillation. Fine-tuning MobileNet-0.25 with distillation reduced the error by 1.1% absolute, and fine-tuning MobileNet-0.5 with distillation reduced the error by 1.3% absolute, compared to fine-tuning without distillation. Fine-tuning MobileNetV2-0.35 with distillation reduced the error by 1.0% absolute, and fine-tuning MobileNetV2-1.0 with distillation reduced the error by 1.0% absolute, compared to fine-tuning without distillation.

For all student-teacher pairs, similarity-preserving distillation obtained lower median error than the spatial attention transfer baseline. Table 5 includes a breakdown of how similarity-preserving distillation compares with spatial attention transfer on a per-split basis. On aggregate, similarity-preserving distillation outperformed spatial attention transfer on 19 out of the 30 MobileNet splits and 24 out of the 30 MobileNetV2 splits. The results suggest that there may be a challenging domain shift in the important image areas for the network to attend. Moreover, while attention transfer summarizes the activation map by summing out the channel dimension, similarity-preserving distillation makes use of the full activation map in computing the similarity-based distillation loss, which may be more robust in the presence of a domain shift in attention.

### 3.3. CINIC-10

The CINIC-10 dataset [5] is designed to be a middle option relative to CIFAR-10 and ImageNet: it is composed of 32x32 images in the style of CIFAR-10, but at a total of 270,000 images its scale is closer to that of ImageNet. We adopted CINIC-10 for rapid experimentation because several GPU-months would have been required to perform full held-out validation and training on ImageNet for our method and all baselines.

For the student and teacher architectures, we experimented with variants of the state-of-the-art mobile architecture ShuffleNetV2 [23]. The ShuffleNetV2 networks are summarized in Table 6. We used the standard training-

Output size	ShuffleNetV2-0.5	ShuffleNetV2-1.0	ShuffleNetV2-2.0
$32 \times 32$	$3 \times 3, 24$	$3 \times 3, 24$	$3 \times 3, 24$
$16 \times 16$	stage( $c = 48, n = 4$ )	stage( $c = 116, n = 4$ )	stage( $c = 244, n = 4$ )
$8 \times 8$	stage( $c = 96, n = 8$ )	stage( $c = 232, n = 8$ )	stage( $c = 488, n = 8$ )
$4 \times 4$	stage( $c = 192, n = 4$ )	stage( $c = 464, n = 4$ )	stage( $c = 976, n = 4$ )
$4 \times 4$	$1 \times 1, 1024$	$1 \times 1, 1024$	$1 \times 1, 2048$
$1 \times 1$	average pool, 10-d fc, softmax		

Table 6. Structure of ShuffleNetV2 networks used in CINIC-10 experiments. The notation ‘stage( $c, n$ )’ denotes a group of ShuffleNetV2 building blocks with  $c$  output channels and  $n$  repeated blocks. Downsampling is performed by strided  $3 \times 3$  depthwise convolutions in the first block of a group.

validation-testing split and set the hyperparameters for similarity-preserving distillation and all baselines by held-out validation (KD:  $\{\alpha = 0.6, T = 16\}$ ; AT:  $\beta = 50$ ; SP:  $\gamma = 2000$ ; KD+SP:  $\{\alpha = 0.6, T = 16, \gamma = 2000\}$ ; AT+SP:  $\{\beta = 30, \gamma = 2000\}$ ). All networks were trained using SGD with Nesterov momentum, a batch size of 96, for 140 epochs with an initial learning rate of 0.01 decayed by a factor of 10 after the 100th and 120th epochs. We applied CIFAR-style data augmentation with horizontal flips and random crops during training.

The results are shown in Table 7. Compared to conventional training with data supervision only, similarity-preserving distillation consistently improved student training outcomes. In particular, training ShuffleNetV2-0.5 with similarity-preserving distillation reduced the error by 1.5% absolute, and training ShuffleNetV2-1.0 with similarity-preserving distillation reduced the error by 1.3% absolute. On an individual basis, all three knowledge distillation approaches achieved comparable results, with a total spread of 0.12% absolute error on ShuffleNetV2-0.5 (for the best results with ShuffleNetV2-1.0 as teacher) and a total spread of 0.06% absolute error on ShuffleNetV2-1.0. However, the lowest error was achieved by combining similarity-preserving distillation with spatial attention transfer. Training ShuffleNetV2-0.5 combining both distillation losses reduced the error by 1.9% absolute, and training ShuffleNetV2-1.0 combining both distillation losses reduced the error by 1.4% absolute. This result shows that similarity-preserving distillation complements attention transfer and captures teacher knowledge that is not fully encoded in spatial attention maps.

**Sensitivity analysis.** Figure 4 illustrates how the performance of similarity-preserving distillation is affected by the choice of hyperparameter  $\gamma$ . We plot the top-1 errors on the CINIC-10 test set for ShuffleNetV2-0.5 and ShuffleNetV2-1.0 students trained with  $\gamma$  ranging from 100 to 16,000. We observed robust performance over a broad range of values for  $\gamma$ . In all experiments, we set  $\gamma$  by held-out validation.

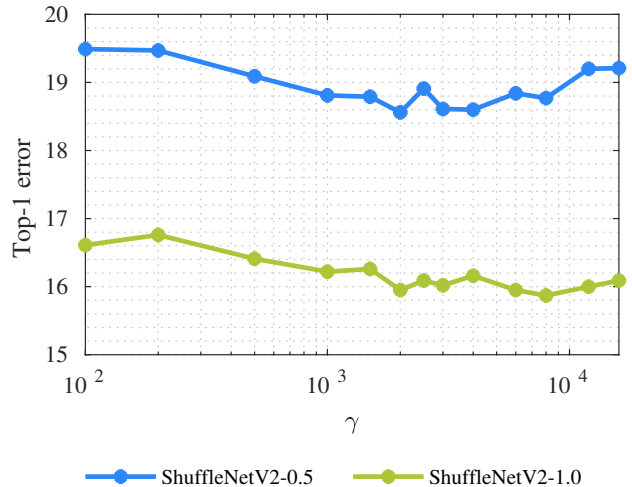


Figure 4. Sensitivity to  $\gamma$  on the CINIC-10 test set for ShuffleNetV2 students.

## 4. Related Work

We presented in this paper a novel distillation loss for capturing and transferring knowledge from a teacher network to a student network. Several prior alternatives [12, 31, 38, 41] are described in the introduction and some key differences are highlighted in Section 2. In addition to the knowledge capture (or loss definition) aspect of distillation studied in this paper, another important open question is the architectural design of students and teachers. In most studies of knowledge distillation, including ours, the student network is a thinner and/or shallower version of the teacher network. Inspired by efficient architectures such as MobileNet and ShuffleNet, Crowley et al. [4] proposed to replace regular convolutions in the teacher network with cheaper grouped and pointwise convolutions in the student. Ashok et al. [1] developed a reinforcement learning approach to learn the student architecture. Polino et al. [28] demonstrated how a quantized student network can be trained using a full-precision teacher network.

There is also innovative orthogonal work exploring al-

Student	Teacher	Student	KD [12]	AT [41]	SP (ours)	KD+SP	AT+SP	Teacher
Sh.NetV2-0.5 (0.4M)	Sh.NetV2-1.0 (1.3M)	20.09	18.62	18.50	18.56	18.35	<b>18.20</b>	17.26
Sh.NetV2-0.5 (0.4M)	Sh.NetV2-2.0 (5.3M)	20.09	18.96	18.78	19.09	18.88	<b>18.43</b>	15.63
Sh.NetV2-1.0 (1.3M)	Sh.NetV2-2.0 (5.3M)	17.26	16.01	15.95	15.95	16.11	<b>15.89</b>	15.63

Table 7. Experiments on CINIC-10 with three different knowledge distillation losses: softened class scores (traditional KD), attention transfer (AT), and similarity preserving (SP). The best result for each experiment is shown in bold. Brackets indicate model size in number of parameters.

ternatives to the usual student-teacher training paradigm. Wang et al. [34] introduced an additional discriminator network, and trained the student, teacher, and discriminator networks together using a combination of distillation and adversarial losses. Lan et al. [18] proposed the on-the-fly native ensemble teacher model, in which the teacher is trained together with multiple students in a multi-branch network architecture. The teacher prediction is a weighted average of the branch predictions.

Knowledge distillation was first introduced as a technique for neural network compression. Resource efficiency considerations have led to a recent increase in interest in efficient neural architectures [13, 14, 23, 32, 43], as well as in algorithms for compressing trained deep networks. Weight pruning methods [11, 20, 22, 24, 33, 35, 39] remove unimportant weights from the network, sparsifying the network connectivity structure. The induced sparsity is unstructured when individual connections are pruned, or structured when entire channels or filters are pruned. Unstructured sparsity usually results in better accuracy but requires specialized sparse matrix multiplication libraries [26] or hardware engines [10] in practice. Quantized networks [8, 15, 17, 30, 42, 45], such as fixed-point, binary, ternary, and arbitrary-bit networks, encode weights and/or activations using a small number of bits, or at lower precision. Fractional or arbitrary-bit quantization [9, 17] encodes individual weights at different precisions, allowing multiple precisions to be used within a single network layer. Low-rank factorization methods [6, 7, 16, 27, 44] produce compact low-rank approximations of filter matrices. Techniques from different categories have also been optimized jointly or combined sequentially to achieve higher compression rates [7, 11, 33].

State-of-the-art network compression methods can achieve significant reductions in network size, in some cases by an order of magnitude, but often require specialized software or hardware support. For example, unstructured pruning requires optimized sparse matrix multiplication routines to realize practical acceleration [26], platform constraint-aware compression [2, 36, 37] requires hardware simulators or empirical measurements, and arbitrary-bit quantization [9, 17] requires specialized hardware. One of the advantages of knowledge distillation is that it is easily implemented in any off-the-shelf deep learning framework

without the need for extra software or hardware. Moreover, distillation can be integrated with other network compression techniques for further gains in performance [28].

## 5. Conclusion

We proposed similarity-preserving knowledge distillation: a novel form of knowledge distillation that aims to preserve pairwise similarities in the student’s representation space, instead of mimicking the teacher’s representation space. Our experiments demonstrate the potential of similarity-preserving distillation in improving the training outcomes of student networks compared to training with only data supervision (e.g. ground truth labels). Moreover, in a transfer learning setting, when traditional class score based distillation is not directly applicable, we have shown that similarity-preserving distillation provides a robust solution to the challenging domain shift problem. We have also shown that similarity-preserving distillation complements the state-of-the-art attention transfer method and captures teacher knowledge that is not fully encoded in spatial attention maps. We believe that similarity-preserving distillation can provide a simple yet effective drop-in replacement for (or complement to) traditional forms of distillation in a variety of application areas, including model compression [28], privileged learning [21], adversarial defense [25], and learning with noisy data [19].

**Future directions.** As future work, we plan to explore similarity-preserving knowledge distillation in semi-supervised and omni-supervised [29] learning settings. Since similarity-preserving distillation does not require labels, it is possible to distill further knowledge from the teacher using auxiliary images without annotations. For example, the supervised loss (e.g. cross-entropy) can be computed using the usual annotated training set, while the distillation loss can be computed using an auxiliary set of unlabelled web images. In this setting, the distillation loss is analogous to the reconstruction or unsupervised loss in semi-supervised learning.

## References

- [1] A. Ashok, N. Rhinehart, F. Beainy, and K. M. Kitani. N2N learning: Network to network compression via policy gradi-



- ent reinforcement learning. In *International Conference on Learning Representations*, 2018.
- [2] C. Chen, F. Tung, N. Vedula, and G. Mori. Constraint-aware deep neural network compression. In *European Conference on Computer Vision*, 2018.
  - [3] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi. Describing textures in the wild. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2014.
  - [4] E. J. Crowley, G. Gray, and A. Storkey. Moonshine: Distilling with cheap convolutions. In *Advances in Neural Information Processing Systems*, 2018.
  - [5] L. N. Darlow, E. J. Crowley, A. Antoniou, and A. J. Storkey. CINIC-10 is not ImageNet or CIFAR-10. arXiv:1810.03505, 2018.
  - [6] E. Denton, W. Zaremba, J. Bruna, Y. LeCun, and R. Fergus. Exploiting linear structure within convolutional networks for efficient evaluation. In *Advances in Neural Information Processing Systems*, 2014.
  - [7] A. Dubey, M. Chatterjee, and N. Ahuja. Coreset-based neural network compression. In *European Conference on Computer Vision*, 2018.
  - [8] J. Faraone, N. Fraser, M. Blott, and P. H. W. Leong. SYQ: Learning symmetric quantization for efficient deep neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
  - [9] J. Fromm, S. Patel, and M. Philipose. Heterogeneous bitwidth binarization in convolutional neural networks. In *Advances in Neural Information Processing Systems*, 2018.
  - [10] S. Han, X. Liu, H. Mao, J. Pu, A. Pedram, M. A. Horowitz, and W. J. Dally. EIE: Efficient inference engine on compressed deep neural network. In *ACM/IEEE International Symposium on Computer Architecture*, 2016.
  - [11] S. Han, H. Mao, and W. J. Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. In *International Conference on Learning Representations*, 2016.
  - [12] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. arXiv:1503.02531, 2015.
  - [13] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam. MobileNets: Efficient convolutional neural networks for mobile vision applications. arXiv:1704.04861, 2017.
  - [14] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer. SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <0.5mb model size. arXiv:1602.07360, 2016.
  - [15] B. Jacob, S. Kligys, B. Chen, M. Zhu, M. Tang, A. Howard, H. Adam, and D. Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
  - [16] M. Jaderberg, A. Vedaldi, and A. Zisserman. Speeding up convolutional neural networks with low rank expansions. In *British Machine Vision Conference*, 2014.
  - [17] S. Khoram and J. Li. Adaptive quantization of neural networks. In *International Conference on Learning Representations*, 2018.
  - [18] X. Lan, X. Zhu, and S. Gong. Knowledge distillation by on-the-fly native ensemble. In *Advances in Neural Information Processing Systems*, 2018.
  - [19] Y. Li, J. Yang, Y. Song, L. Cao, J. Luo, and L.-J. Li. Learning from noisy labels with distillation. In *IEEE International Conference on Computer Vision*, 2017.
  - [20] Z. Liu, J. Xu, X. Peng, and R. Xiong. Frequency-domain dynamic pruning for convolutional neural networks. In *Advances in Neural Information Processing Systems*, 2018.
  - [21] D. Lopez-Paz, L. Bottou, B. Schölkopf, and V. Vapnik. Unifying distillation and privileged information. In *International Conference on Learning Representations*, 2016.
  - [22] J.-H. Luo, J. Wu, and W. Lin. ThiNet: A filter level pruning method for deep neural network compression. In *IEEE International Conference on Computer Vision*, 2017.
  - [23] N. Ma, X. Zhang, H.-T. Zheng, and J. Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In *European Conference on Computer Vision*, 2018.
  - [24] S. Narang, G. Diamos, S. Sengupta, and E. Elsen. Exploring sparsity in recurrent neural networks. In *International Conference on Learning Representations*, 2017.
  - [25] N. Papernot, P. McDaniel, X. Wu, S. Jha, and A. Swami. Distillation as a defense to adversarial perturbations against deep neural networks. In *IEEE Symposium on Security and Privacy*, 2017.
  - [26] J. Park, S. Li, W. Wen, P. Tang, H. Li, Y. Chen, and P. Dubey. Faster CNNs with direct sparse convolutions and guided pruning. In *International Conference on Learning Representations*, 2017.
  - [27] B. Peng, W. Tan, Z. Li, S. Zhang, D. Xie, and S. Pu. Extreme network compression via filter group approximation. In *European Conference on Computer Vision*, 2018.
  - [28] A. Polino, R. Pascanu, and D. Alistarh. Model compression via distillation and quantization. In *International Conference on Learning Representations*, 2018.
  - [29] I. Radosavovic, P. Dollár, R. Girshick, G. Gkioxari, and K. He. Data distillation: towards omni-supervised learning. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
  - [30] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi. XNOR-Net: ImageNet classification using binary convolutional neural networks. In *European Conference on Computer Vision*, 2016.
  - [31] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio. FitNets: hints for thin deep nets. In *International Conference on Learning Representations*, 2015.
  - [32] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen. MobileNetV2: Inverted residuals and linear bottlenecks. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
  - [33] F. Tung and G. Mori. CLIP-Q: Deep network compression learning by in-parallel pruning-quantization. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
  - [34] X. Wang, R. Zhang, Y. Sun, and J. Qi. KDGAN: Knowledge distillation with generative adversarial networks. In *Advances in Neural Information Processing Systems*, 2018.

- [35] W. Wen, C. Wu, Y. Wang, Y. Chen, and H. Li. Learning structured sparsity in deep neural networks. In *Advances in Neural Information Processing Systems*, 2016.
- [36] T.-J. Yang, Y.-H. Chen, and V. Sze. Designing energy-efficient convolutional neural networks using energy-aware pruning. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [37] T.-J. Yang, A. Howard, B. Chen, X. Zhang, A. Go, M. Sandle, V. Sze, and H. Adam. NetAdapt: Platform-aware neural network adaptation for mobile applications. In *European Conference on Computer Vision*, 2018.
- [38] J. Yim, D. Joo, J. Bae, and J. Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [39] R. Yu, A. Li, C.-F. Chen, J.-H. Lai, V. I. Morariu, X. Han, M. Gao, C.-Y. Lin, and L. S. Davis. NISP: Pruning networks using neuron importance score propagation. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [40] S. Zagoruyko and N. Komodakis. Wide residual networks. In *British Machine Vision Conference*, 2016.
- [41] S. Zagoruyko and N. Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *International Conference on Learning Representations*, 2017.
- [42] D. Zhang, J. Yang, D. Ye, and G. Hua. LQ-Nets: Learned quantization for highly accurate and compact deep neural networks. In *European Conference on Computer Vision*, 2018.
- [43] X. Zhang, X. Zhou, M. Lin, and J. Sun. ShuffleNet: An extremely efficient convolutional neural network for mobile devices. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [44] X. Zhang, J. Zou, X. Ming, K. He, and J. Sun. Efficient and accurate approximations of nonlinear convolutional networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2015.
- [45] A. Zhou, A. Yao, K. Wang, and Y. Chen. Explicit loss-error-aware quantization for low-bit deep neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.