

# CONTRASTIVE REPRESENTATION DISTILLATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Often we wish to transfer representational knowledge from one neural network to another. Examples include distilling a large network into a smaller one, transferring knowledge from one sensory modality to a second, or ensembling a collection of models into a single estimator. Knowledge distillation, the standard approach to these problems, minimizes the KL divergence between the probabilistic outputs of a teacher and student network. We demonstrate that this objective ignores important *structural* knowledge of the teacher network. This motivates an alternative objective by which we train a student to capture significantly more *information* in the teacher’s representation of the data. We formulate this objective as contrastive learning. Experiments demonstrate that our resulting new objective outperforms knowledge distillation on a variety of knowledge transfer tasks, including single model compression, ensemble distillation, and cross-modal transfer. When combined with knowledge distillation, our method sets a state of the art in many transfer tasks, sometimes even outperforming the *teacher* network.

## 1 INTRODUCTION

Knowledge distillation (KD) transfers knowledge from one deep learning model (the teacher) to another (the student). The objective originally proposed by Hinton et al. (2015); Buciluă et al. (2006) minimizes the KL divergence between the teacher and student outputs. This formulation makes intuitive sense when the output is a distribution, e.g., a probability mass function over classes. However, often we instead wish to transfer knowledge about a *representation*. For example, in the problem of “cross-modal distillation”, we may wish to transfer the representation of an image processing network to a sound (Aytar et al., 2016) or to depth (Gupta et al., 2015) processing network, such that deep features for an image and the associated sound or depth features are highly correlated. In such cases, the KL divergence is undefined.

Representational knowledge is *structured* – the dimensions exhibit complex interdependencies. The original KD objective introduced in Buciluă et al. (2006); Hinton et al. (2015) treats all dimensions as independent, conditioned on the input. Let  $\mathbf{y}^T$  be the output of the teacher and  $\mathbf{y}^S$  be the output of the student. Then the original KD objective function,  $\psi$ , has the fully factored form:  $\psi(\mathbf{y}^S, \mathbf{y}^T) = \sum_i \phi_i(\mathbf{y}_i^S, \mathbf{y}_i^T)^*$ . Such a factored objective is insufficient for transferring structural knowledge, i.e. dependencies between output dimensions  $i$  and  $j$ . This is similar to the situation in image generation where an  $L_2$  objective produces blurry results, due to independence assumptions between output dimensions.

To overcome this problem, we would like an objective that capture correlations and higher order output dependencies. To achieve this, in this paper we leverage the family of *contrastive* objectives (Gutmann & Hyvärinen, 2010; Oord et al., 2018; Arora et al., 2019; Hjelm et al., 2018). These objective functions have been used successfully in recent years for density estimation and representation learning, especially in self-supervised settings. Here we adapt them to the task of knowledge distillation from one deep network to another. We show that it is important to work in representation space, similar to recent works such as Zagoruyko & Komodakis (2016a); Romero et al. (2014). However, note that the loss functions used in those works do not explicitly try to capture correlations or higher-order dependencies in representational space.

Our objective maximizes a lower-bound to the mutual information between the teacher and student representations. We find that this results in better performance on several knowledge transfer tasks.

\*In particular, in Hinton et al. (2015),  $\phi_i(a, b) = -a \log b \quad \forall i$

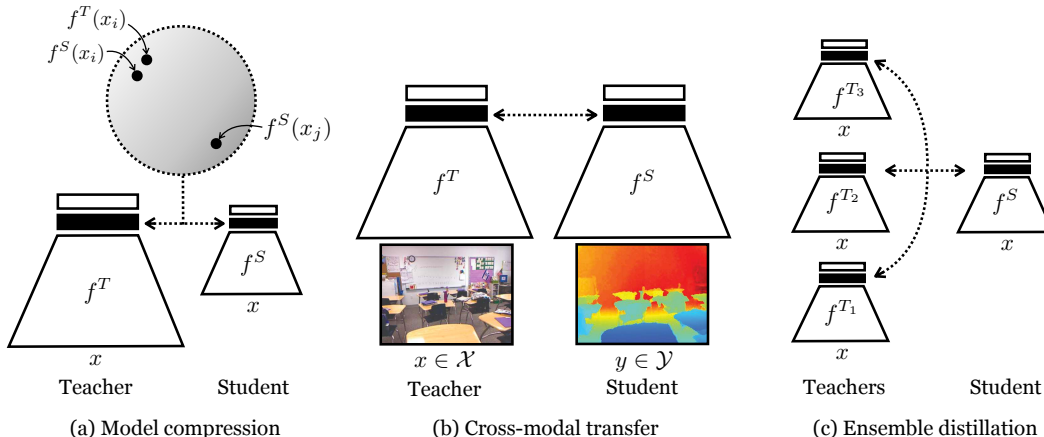


Figure 1: The three distillation settings we consider: (a) compressing a model, (b) transferring knowledge from one modality (e.g., RGB) to another (e.g., depth), (c) distilling an ensemble of nets into a single network. The contrastive objective encourages the teacher and student to map the same input to close representations (in some metric space), and different inputs to distant representations, as indicated in the shaded circle.

We conjecture that this is because the contrastive objective better transfers all the information in the teacher’s representation, rather than only transferring knowledge about conditionally independent output class probabilities. Somewhat surprisingly, the contrastive objective even improves results on the originally proposed task of distilling knowledge about class probabilities, for example, compressing a large CIFAR100 network into a smaller one that performs almost as well. We believe this is because the correlations between different class probabilities contains useful information that regularizes the learning problem. Our paper forges a connection between two literatures that have evolved mostly independently: knowledge distillation and representation learning. This connection allows us to leverage strong methods from representation learning to significantly improve the SOTA on knowledge distillation. Our contributions are:

1. We propose a contrastive based objective for transfer of knowledge between one or more teacher networks to a student network and show that it captures representation structure much more effectively (see Fig. 2). We formulate different variants of our loss function for applications in model compression, cross-modal transfer, and compressing ensembles.
2. We conduct extensive experiments with 12 recent state-of-the-art distillation objectives on various student-teacher combinations. Surprisingly, we found that original Knowledge Distillation (KD) Hinton et al. (2015) is still a strong baseline and none of the recent distillation objectives can consistently outperform KD (without combining with it). Meanwhile, our method consistently surpasses KD and other methods.

## 2 RELATED WORK

The seminal work of Buciluă et al. (2006) and Hinton et al. (2015) introduced the idea of knowledge distillation between large, cumbersome models into smaller, faster models without losing too much generalization power. The general motivation was that at training time, the availability of computation allows “slop” in model size, and potentially faster learning. But computation and memory constraints at inference time necessitate the use of smaller models. Buciluă et al. (2006) achieve this by matching output logits; Hinton et al. (2015) introduced the idea of temperature in the softmax outputs to better represent smaller probabilities in the output of a single sample. These smaller probabilities provide useful information about the learned representation of the teacher model; some tradeoff between large temperatures (which increase entropy) or small temperatures tend to provide the highest transfer of knowledge between student and teacher. The method in Li et al. (2014) was also closely related to Hinton et al. (2015).

Attention transfer (Zagoruyko & Komodakis, 2016a) focuses on the features maps of the network as opposed to the output logits. Here the idea is to elicit similar response patterns in the teacher and student feature maps (called “attention”). However, only feature maps with the same spatial resolution can be combined in this approach, which is a significant limitation since it requires student and teacher networks with very similar architectures. This technique achieves state of the art results

for distillation (as measured by the generalization of the student network). FitNets (Romero et al., 2014) also deal with intermediate representations by using regressions to guide the feature activations of the student network. Since Zagoruyko & Komodakis (2016a) do a weighted form of this regression, they tend to perform better. Other papers (Yim et al., 2017; Huang & Wang, 2017; Kim et al., 2018; Yim et al., 2017; Huang & Wang, 2017; Ahn et al., 2019) have enforced various criteria based on representations. The contrastive objective we use in this paper is related to the InfoNCE and NCE objectives introduced in (Oord et al., 2018; Gutmann & Hyvärinen, 2010). (Oord et al., 2018) use contrastive learning in the context of self-supervised learning of representations. They show that their objective maximizes a lower bound on mutual information. A very related approach is used in Hjelm et al. (2018). InfoNCE and NCE are closely related but distinct from adversarial learning (Goodfellow et al., 2014). In Goodfellow (2014), it is shown that the NCE objective of Gutmann & Hyvärinen (2010) can lead to maximum likelihood learning, but not the adversarial objective.

### 3 METHOD

The key idea of contrastive learning is very general: learn a representation that is close in some metric space for “positive” pairs and push apart the representation between “negative” pairs. Fig. 1 gives a visual explanation for how we structure contrastive learning for the three tasks we consider: model compression, cross-modal transfer and ensemble distillation.

#### 3.1 CONTRASTIVE LOSS

Given two deep neural networks, a teacher  $f^T$  and a student  $f^S$ . Let  $x$  be the network input; we denote representations at the penultimate layer (before logits) as  $f^T(x)$  and  $f^S(x)$  respectively. Let  $x_i$  represent a training sample, and  $x_j$  another randomly chosen sample. We would like to push closer the representations  $f^S(x_i)$  and  $f^T(x_i)$  while pushing apart  $f^S(x_i)$  and  $f^T(x_j)$ . For ease of notation, we define random variables  $S$  and  $T$  for the student and teacher’s representations of the data respectively:

$$x \sim p_{\text{data}}(x) \quad \triangleleft \quad \mathbf{data} \quad (1)$$

$$S = f^S(x) \quad \triangleleft \quad \mathbf{student's representation} \quad (2)$$

$$T = f^T(x) \quad \triangleleft \quad \mathbf{teacher's representation} \quad (3)$$

Intuitively speaking, we will consider the joint distribution  $p(S, T)$  and the product of marginal distributions  $p(S)p(T)$ , so that, by maximizing KL divergence between these distributions, we can maximize the *mutual information* between student and teacher representations. To setup an appropriate loss that can achieve this aim, let us define a distribution  $q$  with latent variable  $C$  which decides whether a tuple  $(f^T(x_i), f^S(x_j))$  was drawn from the joint ( $C = 1$ ) or product of marginals ( $C = 0$ ):

$$q(T, S|C = 1) = p(T, S), \quad q(T, S|C = 0) = p(T)p(S) \quad (4)$$

Now, suppose in our data, we are given 1 congruent pair (drawn from the joint distribution, i.e. the same input provided to  $T$  and  $S$ ) for every  $N$  incongruent pairs (drawn from the product of marginals; independent randomly drawn inputs provided to  $T$  and  $S$ ). Then the priors on the latent  $C$  are:

$$q(C = 1) = \frac{1}{N + 1}, \quad q(C = 0) = \frac{N}{N + 1} \quad (5)$$

By simple manipulation and Bayes’ rule, the posterior for class  $C = 1$  is given by:

$$q(C = 1|T, S) = \frac{q(T, S|C = 1)q(C = 1)}{q(T, S|C = 0)q(C = 0) + q(T, S|C = 1)q(C = 1)} \quad (6)$$

$$= \frac{p(T, S)}{p(T, S) + Np(T)p(S)} \quad (7)$$

Next, we observe a connection to mutual information as follows:

$$\begin{aligned} \log q(C = 1|T, S) &= \log \frac{p(T, S)}{p(T, S) + Np(T)p(S)} \\ &= -\log\left(1 + N\frac{p(T)p(S)}{p(T, S)}\right) \leq -\log(N) + \log \frac{p(T, S)}{p(T)p(S)} \end{aligned} \quad (8)$$

Then taking expectation on both sides w.r.t.  $p(T, S)$  (equivalently w.r.t.  $q(T, S|C = 1)$ ) and rearranging, gives us:

$$I(T; S) \geq \log(N) + \mathbb{E}_{q(T, S|C=1)} \log q(C = 1|T, S) \quad \triangleleft \quad \mathbf{MI \ bound} \quad (9)$$

where  $I(T; S)$  is the mutual information between the distributions of the teacher and student embeddings. Thus maximizing  $\mathbb{E}_{q(T, S|C=1)} \log q(C = 1|T, S)$  w.r.t. the parameters of the student network  $S$  increases a lower bound on mutual information. However, we do not know the true distribution  $q(C = 1|T, S)$ ; instead we estimate it by fitting a model  $h : \{\mathcal{T}, \mathcal{S}\} \rightarrow [0, 1]$  to samples from the data distribution  $q(C = 1|T, S)$ , where  $\mathcal{T}$  and  $\mathcal{S}$  represent the domains of the embeddings. We maximize the log likelihood of the data under this model (a binary classification problem):

$$\mathcal{L}_{critic}(h) = \mathbb{E}_{q(T, S|C=1)}[\log h(T, S)] + N\mathbb{E}_{q(T, S|C=0)}[1 - \log(h(T, S))] \quad (10)$$

$$h^* = \arg \max_h \mathcal{L}_{critic}(h) \quad \triangleleft \quad \mathbf{optimal \ critic} \quad (11)$$

We term  $h$  the *critic* since we will be learning representations that optimize the critic’s score. Assuming sufficiently expressive  $h$ ,  $h^*(T, S) = q(C = 1|T, S)$  (via Gibb’s inequality; see Sec. 6.2.1 for proof), so we can rewrite Eq. 9 in terms of  $h^*$ :

$$I(T; S) \geq \log(N) + \mathbb{E}_{q(T, S|C=1)}[\log h^*(T, S)] \quad (12)$$

Therefore, we see that the optimal critic is an estimator whose expectation lower-bounds mutual information. We wish to learn a student that maximizes the mutual information between its representation and the teacher’s, suggesting the following optimization problem:

$$f^{S*} = \arg \max_{f^S} \mathbb{E}_{q(T, S|C=1)}[\log h^*(T, S)] \quad (13)$$

An apparent difficulty here is that the optimal critic  $h^*$  depends on the current student. We can circumvent this difficulty by weakening the bound in (12) to:

$$I(T; S) \geq \log(N) + \mathbb{E}_{q(T, S|C=1)}[\log h^*(T, S)] + N\mathbb{E}_{q(T, S|C=0)}[\log(1 - h^*(T, S))] \quad (14)$$

$$= \log(N) + \mathcal{L}_{critic}(h^*) = \log(N) + \max_h \mathcal{L}_{critic}(h) \quad (15)$$

$$\geq \log(N) + \mathcal{L}_{critic}(h) \quad (16)$$

The first line comes about by simply adding  $N\mathbb{E}_{q(T, S|C=0)}[\log(1 - h^*(T, S))]$  to the bound in (12). This term is strictly negative, so the inequality holds. The last line follows from the fact that  $\mathcal{L}_{critic}(h^*)$  upper-bounds  $\mathcal{L}_{critic}(h)$ . Optimizing (15) w.r.t. the student we have:

$$f^{S*} = \arg \max_{f^S} \max_h \mathcal{L}_{critic}(h) \quad \triangleleft \quad \mathbf{our \ final \ learning \ problem} \quad (17)$$

$$= \arg \max_{f^S} \max_h \mathbb{E}_{q(T, S|C=1)}[\log h(T, S)] + N\mathbb{E}_{q(T, S|C=0)}[\log(1 - h(T, S))] \quad (18)$$

which demonstrates that we may jointly optimize  $f^S$  at the same time as we learn  $h$ . We note that due to (16),  $f^{S*} = \arg \max_{f^S} \mathcal{L}_{critic}(h)$ , for any  $h$ , *also* is a representation that optimizes a lower-bound (a weaker one) on mutual information, so our formulation does not rely on  $h$  being optimized perfectly.

We may choose to represent  $h$  with any family of functions that satisfy  $h : \{\mathcal{T}, \mathcal{S}\} \rightarrow [0, 1]$ . In practice, we use the following:

$$h(T, S) = \frac{e^{g^T(T)'g^S(S)/\tau}}{e^{g^T(T)'g^S(S)/\tau} + \frac{N}{M}} \quad (19)$$

where  $M$  is the cardinality of the dataset and  $\tau$  is a temperature that adjusts the concentration level. In practice, since the dimensionality of  $S$  and  $T$  may be different,  $g^S$  and  $g^T$  linearly transform them into the same dimension and further normalize them by  $\mathcal{L}$ -2 norm before the inner product. space The form of Eq. (18) is inspired by NCE (Gutmann & Hyvärinen, 2010). Our formulation is similar to the InfoNCE loss (Oord et al., 2018) in that we maximize a lower bound on the mutual information. However we use a different objective and bound, which in preliminary experiments we found to be more effective than InfoNCE.

**Implementation.** Theoretically, larger  $N$  in Eq. 16 leads to tighter lower bound on MI. In practice, to avoid very large batch size, we implement a memory buffer that stores latent features of each data sample computed from previous batches. Therefore, during training we can efficiently retrieve a large number of negative samples from the memory buffer without recomputing their features.

### 3.2 KNOWLEDGE DISTILLATION OBJECTIVE

The knowledge distillation loss was proposed in Hinton et al. (2015). In addition to the regular cross-entropy loss between the student output  $y^S$  and one-hot label  $y$ , it asks the student network output to be as similar as possible to the teacher output by minimizing the cross-entropy between their output probabilities. The complete objective is:

$$\mathcal{L}_{KD} = (1 - \alpha)H(y, y^S) + \alpha\rho^2 H(\sigma(z^T/\rho), \sigma(z^S/\rho)) \quad (20)$$

where  $\rho$  is the temperature,  $\alpha$  is a balancing weight,  $\sigma$  is softmax function.

### 3.3 CROSS-MODAL TRANSFER LOSS

In the cross-modal transfer task shown in Fig. 1(b), a teacher network is trained on a source modality  $\mathcal{X}$  with large-scale labeled dataset. We then wish to transfer the knowledge to a student network, but adapt it to another dataset or modality  $\mathcal{Y}$ . But the features of the teacher network are still valuable to help with learning of the student on another domain. In this transfer task, we use the contrastive loss Eq. 10 to match the features of the student and teacher. Additionally, we also consider other distillation objectives discussed in previous section. Such transfer is conducted on a paired but unlabeled dataset  $D = \{(x_i, y_i) | i = 1, \dots, L, x_i \in \mathcal{X}, y_i \in \mathcal{Y}\}$ . In this scenario, there is no true label  $y$  of such data for the original training task on the source modality, and therefore we ignore the  $H(y, y^S)$  term in all objectives that we test. Prior cross-modal work Aytar et al. (2016); Hoffman et al. (2016b;a) uses either  $L_2$  regression or KL-divergence.

### 3.4 ENSEMBLE DISTILLATION LOSS

In the case of ensemble distillation shown in 1(c), we have  $M > 1$  teacher networks,  $f^{T_i}$  and one student network  $f^S$ . We adopt the contrastive framework by defining multiple pairwise contrastive losses between features of each teacher network  $f^{T_i}$  and the student network  $f^S$ . These losses are summed together to give the final loss (to be minimized):  $\mathcal{L}_{CRD-EN} = H(y, y^S) - \beta \sum_i \mathcal{L}_{critic}(T_i, S)$ .

## 4 EXPERIMENTS

We evaluate our contrastive knowledge distillation (CKD) framework in three knowledge distillation tasks: (a) model compression of a large network to a smaller one; (b) cross-modal knowledge transfer; (c) ensemble distillation from a group of teachers to a single student network.

**Datasets** (1) *CIFAR-100* (Krizhevsky & Hinton, 2009) contains 50K training images with 0.5K images per class and 10K test images. (2) *ImageNet* (Deng et al., 2009) provides 1.2 million images from 1K classes for training and 50K for validation. (3) *STL-10* (Coates et al., 2011) consists of a training set of 5K labeled images from 10 classes and 100K unlabeled images, and a test set of 8K images. (4) *TinyImageNet* (Deng et al., 2009) has 200 classes, each with 500 training images and 50 validation images. (5) *NYU-Depth V2* (Silberman et al., 2012) consists of 1449 indoor images, each labeled with dense depth image and semantic map.

### 4.1 MODEL COMPRESSION

**Setup** We experiment on CIFAR-100 and ImageNet with student-teacher combinations of various capacity, such as ResNet (He et al., 2016) or Wide ResNet (WRN) (Zagoruyko & Komodakis, 2016b).

**Results on CIFAR100** Table 1 and Table 2 compare top-1 *accuracies* of different distillation objectives (for details, see Section 6.1). Table 1 investigates students and teachers of the same architectural style, while Table 2 focuses on students and teachers from different architectures. We observe that our loss, which we call CRD (Contrastive Representation Distillation), consistently outperforms all other distillation objectives, including the original knowledge distillation (KD). Surprisingly, KD works pretty well and none of the other methods consistently outperform KD on their own. Another observation is that, while switching the teacher student combinations from similar architectures to different architectures, methods that distill intermediate representations tend to perform worse than methods that distill from the last several layers. For example, the Attention Transfer (AT) method performs the worst with MobileNetV2 as the student network, and it even underperforms the vanilla student. In contrast, PKT, SP and CRD that operate on last several layers performs well.

Teacher	WRN-40-2	WRN-40-2	resnet56	resnet110	resnet110	resnet32x4	vgg13
Student	WRN-16-2	WRN-40-1	resnet20	resnet20	resnet32	resnet8x4	vgg8
Teacher	75.61	75.61	72.34	74.31	74.31	79.42	74.64
Student	73.26	71.98	69.06	69.06	71.14	72.50	70.36
KD*	74.92	73.54	70.66	70.67	73.08	73.33	72.98
FitNet*	73.58 (↓)	72.24 (↓)	69.21 (↓)	68.99 (↓)	71.06 (↓)	73.50 (↑)	71.02 (↓)
AT	74.08 (↓)	72.77 (↓)	70.55 (↓)	70.22 (↓)	72.31 (↓)	73.44 (↑)	71.43 (↓)
SP	73.83 (↓)	72.43 (↓)	69.67 (↓)	70.04 (↓)	72.69 (↓)	72.94 (↓)	72.68 (↓)
CC	73.56 (↓)	72.21 (↓)	69.63 (↓)	69.48 (↓)	71.48 (↓)	72.97 (↓)	70.71 (↓)
VID	74.11 (↓)	73.30 (↓)	70.38 (↓)	70.16 (↓)	72.61 (↓)	73.09 (↓)	71.23 (↓)
RKD	73.35 (↓)	72.22 (↓)	69.61 (↓)	69.25 (↓)	71.82 (↓)	71.90 (↓)	71.48 (↓)
PKT	74.54 (↓)	73.45 (↓)	70.34 (↓)	70.25 (↓)	72.61 (↓)	73.64 (↑)	72.88 (↓)
AB	72.50 (↓)	72.38 (↓)	69.47 (↓)	69.53 (↓)	70.98 (↓)	73.17 (↓)	70.94 (↓)
FT*	73.25 (↓)	71.59 (↓)	69.84 (↓)	70.22 (↓)	72.37 (↓)	72.86 (↓)	70.58 (↓)
FSP*	72.91 (↓)	0.00 (↓)	69.95 (↓)	70.11 (↓)	71.89 (↓)	72.62 (↓)	70.23 (↓)
NST*	73.68 (↓)	72.24 (↓)	69.60 (↓)	69.53 (↓)	71.96 (↓)	73.30 (↓)	71.53 (↓)
CRD	<b>75.48</b> (↑)	<b>74.14</b> (↑)	<b>71.16</b> (↑)	<b>71.46</b> (↑)	<b>73.48</b> (↑)	<b>75.51</b> (↑)	<b>73.94</b> (↑)

Table 1: Test *accuracy* (%) of student networks on CIFAR100 of a number of distillation methods (ours is CRD); see Appendix for citations of other methods. ↑ denotes outperformance over KD and ↓ denotes underperformance. We note that CRD is the *only* method to always outperform KD (and also outperforms all other methods). We denote by \* methods where we used our reimplementations based on the paper; for all other methods we used author-provided or author-verified code. Average over 5 runs.

Teacher	vgg13	ResNet50	ResNet50	resnet32x4	resnet32x4	WRN-40-2
Student	MobileNetV2	MobileNetV2	vgg8	ShuffleNetV1	ShuffleNetV2	ShuffleNetV1
Teacher	74.64	79.34	79.34	79.42	79.42	75.61
Student	64.6	64.6	70.36	70.5	71.82	70.5
KD*	67.37	67.35	73.81	74.07	74.45	74.83
FitNet*	64.14 (↓)	63.16 (↓)	70.69 (↓)	73.59 (↓)	73.54 (↓)	73.73 (↓)
AT	59.40 (↓)	58.58 (↓)	71.84 (↓)	71.73 (↓)	72.73 (↓)	73.32 (↓)
SP	66.30 (↓)	68.08 (↑)	73.34 (↓)	73.48 (↓)	74.56 (↑)	74.52 (↓)
CC	64.86 (↓)	65.43 (↓)	70.25 (↓)	71.14 (↓)	71.29 (↓)	71.38 (↓)
VID	65.56 (↓)	67.57 (↑)	70.30 (↓)	73.38 (↓)	73.40 (↓)	73.61 (↓)
RKD	64.52 (↓)	64.43 (↓)	71.50 (↓)	72.28 (↓)	73.21 (↓)	72.21 (↓)
PKT	67.13 (↓)	66.52 (↓)	73.01 (↓)	74.10 (↑)	74.69 (↑)	73.89 (↓)
AB	66.06 (↓)	67.20 (↓)	70.65 (↓)	73.55 (↓)	74.31 (↓)	73.34 (↓)
FT*	61.78 (↓)	60.99 (↓)	70.29 (↓)	71.75 (↓)	72.50 (↓)	72.03 (↓)
NST*	58.16 (↓)	64.96 (↓)	71.28 (↓)	74.12 (↑)	74.68 (↑)	74.89 (↑)
CRD	<b>69.73</b> (↑)	<b>69.11</b> (↑)	<b>74.30</b> (↑)	<b>75.11</b> (↑)	<b>75.65</b> (↑)	<b>76.05</b> (↑)

Table 2: Top-1 test *accuracy* (%) of student networks on CIFAR100 of a number of distillation methods (ours is CRD) for transfer across very different teacher and student architectures. CRD outperforms KD and all other methods. Importantly, some methods that require very similar student and teacher architectures perform quite poorly. E.g. FSP (Yim et al., 2017) cannot even be applied; AT (Ba & Caruana, 2014) and FitNet (Zagoruyko & Komodakis, 2016a) perform very poorly etc. We denote by \* methods where we used our reimplementations based on the paper; for all other methods we used author-provided or author-verified code. Average over 3 runs.

In Fig. 2, we compute the difference of cross-correlations between the teacher and student logits; for three different students: randomly initialized, trained by AT, KD or CRD (our method). It is clear that the CRD objective captures the most correlation structure in the logit as shown by the smaller differences between teacher and student. This is reflected in reduced error rates.

**Results on ImageNet** For a fair comparison with Zagoruyko & Komodakis (2016a) and Lan et al. (2018), we adopt the models from these papers, ResNet-34 as the teacher and ResNet-18 as the student. As shown in Table 3, the gap of top-1 accuracy between the teacher and student is 3.56%. The AT

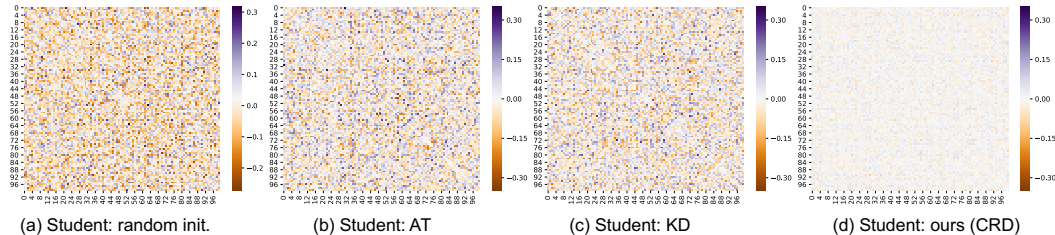


Figure 2: The cross-correlations between class logits of a teacher network are ignored by regular cross-entropy. Distillation frameworks use “soft targets” (Hinton et al., 2015) which effectively capture such correlations and transfer them to the student network, leading to the success of distillation. We visualize here the *difference* of normalized cross-correlation matrices of student and teacher logits, for different student networks on a CIFAR-100 knowledge distillation task: (a) Student trained from random initialization, showing that the teacher and student cross-correlations are very different; (b) Student distilled by attention transfer Zagoruyko & Komodakis (2016a); showing reduced difference (see axis); (c) Student distilled by KL divergence Hinton et al. (2015), also showing reduced difference; (d) Student distilled by our contrastive objective, showing significant matching between student and teacher correlations.

	Teacher	Student	AT	KD	SP	CC	Online KD *	CRD
Top-1	26.69	30.25	29.30	29.34	29.38	30.04	29.45	<b>28.83</b>
Top-5	8.58	10.93	10.00	10.12	10.20	10.83	10.41	<b>9.87</b>

Table 3: Top-1 and Top-5 error rates (%) of student network ResNet-18 on ImageNet validation set. We use ResNet-34 released by PyTorch team as our teacher network, and follow the standard training practice of ImageNet on PyTorch except that we train for 10 more epochs. We compare our CRD with KD Hinton et al. (2015), AT Zagoruyko & Komodakis (2016a) and Online-KD Lan et al. (2018). “\*” reported by the original paper Lan et al. (2018) using an ensemble of ResNets as teacher.

method reduces this gap by 0.95%, while ours narrow it by 1.42%, a 50% relative improvement. Results on ImageNet validates the scalability of our CRD.

	Student	KD	AT	FitNet	CRD	CRD+KD	Teacher
CIFAR100→STL-10	69.7	70.9	70.7	70.3	71.6	<b>72.2</b>	68.6
CIFAR100→TinyImageNet	33.7	33.9	34.2	33.5	<b>35.6</b>	35.5	31.5

Table 4: We transfer the representation learned from CIFAR100 to STL-10 and TinyImageNet datasets by freezing the network and training a linear classifier on top of the last feature layer to perform 10-way (STL-10) or 200-way (TinyImageNet) classification. For this experiment, we use the combination of teacher network WRN-40-2 and student network WRN-16-2. All numbers report accuracies on the target dataset (STL-10 or TinyImageNet).

**Transferability of representations** We are interested in *representations*, and a primary goal of representation learning is to acquire *general* knowledge, that is, knowledge that transfers to tasks or datasets that were unseen during training. Therefore, we test if the representations we distill transfer well. A WRN-16-2 student either distills from a WRN-40-2 teacher, or is trained from scratch on CIFAR100. The student serves as a frozen representation extractor (the layer prior to the logit) for images from STL-10 or TinyImageNet (all images downsampled to 32x32). We then train a *linear* classifier to perform 10-way (for STL-10) or 200-way (for TinyImageNet) classification to quantify the transferability of the representations. We compare CRD with multiple baselines such as KD and AT in Table 3. In general, all distillation methods except FitNet improve transferability of the learned representations on both STL-10 and TinyImageNet. While the teacher performs the best on the original CIFAR100 dataset, its representations transfer the worst to the other two datasets. This is perhaps the teacher’s representations are biased towards the original task. Surprisingly, the student with CRD+KD distillation not only matches its teacher on CIFAR100 (see Table 4), but also transfers much better than the teacher, e.g., 3.6% improvement (STL-10) and 4.1% on TinyImageNet.

## 4.2 CROSS-MODAL TRANSFER

We consider a practical setting where modality  $\mathcal{X}$  has large amount of labeled data while modality  $\mathcal{Y}$  does not. Transferring knowledge from  $\mathcal{X}$  to  $\mathcal{Y}$  is a common challenge. For example, while large-scale RGB datasets are easily accessible, other modalities such as depth images are much harder to label at

Metric (%)	Random Init.	KD	KD+AT	FitNet	CRD
Pix. Acc.	56.4	58.9	60.1	60.8	<b>61.6</b>
mIoU	35.8	38.0	39.5	40.7	<b>41.8</b>

Table 5: Performance on the task of using depth to predict semantic segmentation labels. We initialize the depth network either randomly or by distilling from a ImageNet pre-trained ResNet-18 teacher.

scale but have wide applications. We demonstrate the potential of CRD for cross-modal transfer in two scenarios: (a) transfer from luminance to chrominance; (b) transfer from RGB to depth images.

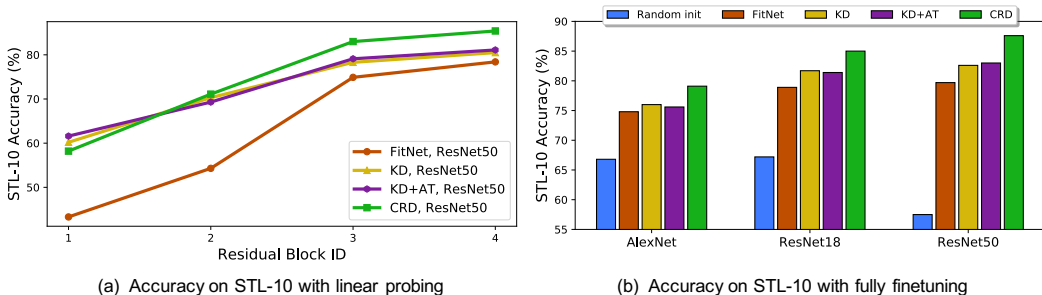


Figure 3: Top-1 classification accuracy on STL-10 using *chrominance* image (*ab* channel in *Lab* color space). We initialize the *chrominance* network randomly or by distilling from a *luminance* network, trained with large-scale labeled images. We evaluate distillation performance by (a) linear probing and (b) fully finetuning.

**Transferring from Luminance to Chrominance.** We work on *Lab* color space, where *L* represents Luminance and *ab* Chrominance. We first train an *L* network on TinyImageNet with supervision. Then we transfer knowledge from this *L* network to *ab* network on the unlabeled set of STL-10 with different objectives, including FitNet, KD, KD+AT and CRD. For convenience, we use the same architecture for student and teacher (they can also be different). Finally, we evaluate the knowledge of *ab* network by two means: (1) *linear probing*: we freeze the *ab* network and train a linear classifier on top of features from different layers to perform 10-way classification on STL-10 *ab* images. This is a common practice (Alain & Bengio, 2016; Zhang et al., 2017) to evaluate the quality of network representations; (b) *fully finetuning*: we fully finetune the *ab* network to obtain the best accuracy. We also use as a baseline the *ab* network that is randomly initialized rather than distilled. Architectures investigated include AlexNet, ResNet-18 and ResNet-50. The results shown in Figure 3 show that CRD is more efficient for transferring inter-modal knowledge than other methods. Besides, we also note KD+AT does not improve upon KD, possibly because attention of luminance and chrominance are different and harder to transfer.

**Transferring from RGB to Depth.** We transfer the knowledge of a ResNet-18 teacher pretrained on ImageNet to a 5-layer student CNN operating on depth images. We follow a similar transferring procedure on NYU-Depth training set, except that we use a trick of contrasting between local and global features, proposed by Hjelm et al. (2018), to overcome the problem of insufficient data samples in the depth domain. Then the student network is further trained to predict semantic segmentation maps from depth images. We note that both knowledge transfer and downstream training are conducted on the same set of images, i.e., the training set. Table 5 reports the average pixel prediction accuracy and mean Intersection-over-Union across all classes. All distillation methods can transfer knowledge from the RGB ResNet to the depth CNN. FitNet surpasses KD and KD+AT. CRD significantly outperforms all other methods.

#### 4.3 DISTILLATION FROM AN ENSEMBLE

Better classification performance is often achieved by ensembles of deep networks, but these are usually too expensive for inference time, and distillation into a single network is a desirable task. We investigate the KL-divergence based KD and our CRD for this task, using the loss of Sec. 3.4. The network structures of each teacher and student are identical here, but an ensemble of multiple teachers can still provide rich knowledge for the student. To compare between KD and CRD on CIFAR100 dataset, we use WRN-16-2 and ResNet-20, whose single model error rates are 26.7% and 30.9% respectively. The results of distillation are presented in Figure 4, where we vary the number of ensembled teachers. CRD with 8 teachers decreases the error rate of WRN-16-2 to 23.7% and ResNet20 to 28.3%. In addition, CRD works consistently better than KD in all settings we test. These



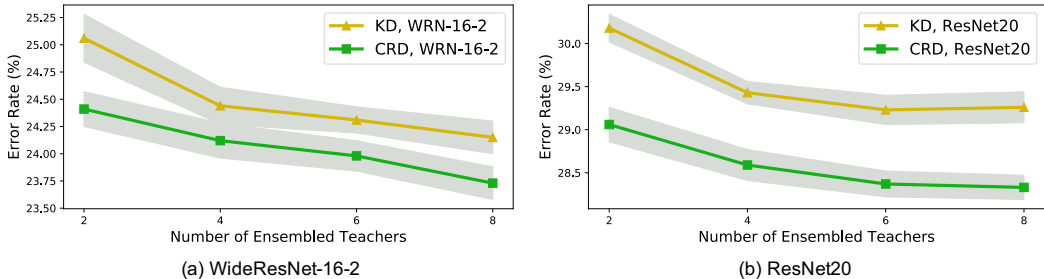


Figure 4: Distillation from an ensemble of teachers. We vary the number of ensembled teachers and compare KD with our CRD by using (a) WRN-16-2 and (b) ResNet20. Our CRD consistently achieves lower error rate. Variations across multiple runs are also shown.

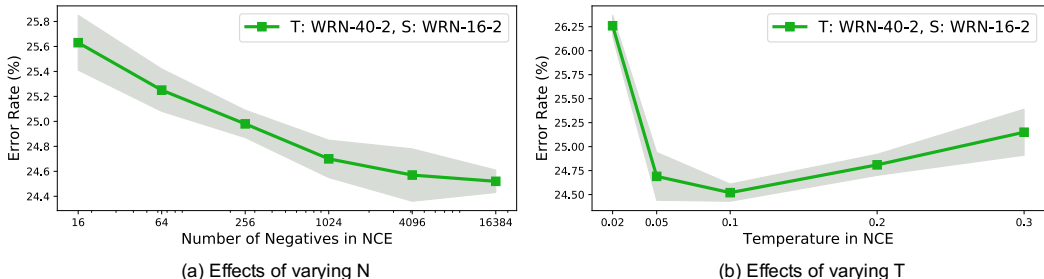


Figure 5: Effects of varying the number of negatives, shown in (a), or the temperature, shown in (b).

observations suggest that CRD is capable of distilling an ensemble of models into a single one which performs significantly better than a model of the same size that is trained from scratch.

#### 4.4 HYPER-PARAMETERS AND COMPUTATION OVERHEAD

There are two main hyper-parameters in our contrastive objectives: (1) the number of negative samples  $N$  in Eq. 18, and (2) temperature  $\tau$  which modulates the softmax probability. We adopt WRN-40-2 as teacher and WRN-16-2 as student for parameter analysis. Experiments are conducted on CIFAR100 and the results are shown in Figure 5.

**Number of negatives  $N$**  We have validated different  $N$ : 16, 64, 256, 1024, 4096, 16384. As shown in Figure 5(a), increasing  $N$  leads to improved performance. However, the difference of error rate between  $N = 4096$  and  $N = 16384$  is less than 0.1%. Therefore, we use  $N = 16384$  for reporting the accuracy while in practice  $N = 4096$  should suffice.

**Temperature  $\tau$**  We varied  $\tau$  between 0.02 and 0.3. As Figure 5(b) illustrates, both extremely high or low temperature lead to a sub-optimal solution. In general, temperatures between 0.05 and 0.2 work well on CIFAR100. All experiments but those on ImageNet use a temperature of 0.1. For ImageNet, we use  $\tau = 0.07$ . The optimal temperature may vary across different datasets and require further tuning.

**Computational Cost** We use ResNet-18 on ImageNet for illustration. CRD uses extra 260 MFLOPs, which is about 12% of the original 2 GFLOPs. In practice, we did not notice significant difference of training time on ImageNet (e.g., 1.75 epochs/hour v.s. 1.67 epochs/hour on two Titan-V GPUs). The memory bank for storing all 128-d features of ImageNet only costs around 600MB memory, and therefore we store it on GPU memory.

## 5 CONCLUSION

We have developed a novel technique for neural network distillation, using the concept of contrastive objectives, which are usually used for representation learning. We experimented with our objective on a number of applications such as model compression, cross-modal transfer and ensemble distillation, outperforming other distillation objectives by significant margins in all these tasks. Our contrastive objective is the only distillation objective that consistently outperforms knowledge distillation across a wide variety of knowledge transfer tasks. Prior objectives only surpass KD *when combined* with KD. Contrastive learning is a simple and effective objective with practical benefits.

## REFERENCES

- Sungsoo Ahn, Shell Xu Hu, Andreas Damianou, Neil D Lawrence, and Zhenwen Dai. Variational information distillation for knowledge transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 9163–9171, 2019. 3, 12, 14
- Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes. *arXiv preprint arXiv:1610.01644*, 2016. 8
- Sanjeev Arora, Hrishikesh Khandeparkar, Mikhail Khodak, Orestis Plevrakis, and Nikunj Saunshi. A theoretical analysis of contrastive unsupervised representation learning. *arXiv preprint arXiv:1902.09229*, 2019. 1
- Yusuf Aytar, Carl Vondrick, and Antonio Torralba. Soundnet: Learning sound representations from unlabeled video. In *Advances in neural information processing systems*, 2016. 1, 5
- Jimmy Ba and Rich Caruana. Do deep nets really need to be deep? In *Advances in neural information processing systems*, pp. 2654–2662, 2014. 6, 12, 13
- Cristian Buciluă, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 535–541. ACM, 2006. 1, 2
- Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pp. 215–223, 2011. 5
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009. 5
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pp. 2672–2680, 2014. 3
- Ian J Goodfellow. On distinguishability criteria for estimating generative models. *arXiv preprint arXiv:1412.6515*, 2014. 3
- Saurabh Gupta, Judy Hoffman, and Jitendra Malik. Cross modal distillation for supervision transfer. *CoRR*, abs/1507.00448, 2015. URL <http://arxiv.org/abs/1507.00448>. 1
- Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pp. 297–304, 2010. 1, 3, 4
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016. 5, 13
- Byeongho Heo, Minsik Lee, Sangdoon Yun, and Jin Young Choi. Knowledge transfer via distillation of activation boundaries formed by hidden neurons. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 3779–3787, 2019. 12, 14
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 1, 2, 5, 7, 12, 14, 15
- R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. *arXiv preprint arXiv:1808.06670*, 2018. 1, 3, 8
- Judy Hoffman, Saurabh Gupta, and Trevor Darrell. Learning with side information through modality hallucination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 826–834, 2016a. 5

- Judy Hoffman, Saurabh Gupta, Jian Leong, Sergio Guadarrama, and Trevor Darrell. Cross-modal adaptation for rgb-d detection. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 5032–5039. IEEE, 2016b. [5](#)
- Zehao Huang and Naiyan Wang. Like what you like: Knowledge distill via neuron selectivity transfer. *arXiv preprint arXiv:1707.01219*, 2017. [3](#), [12](#), [14](#)
- Jangho Kim, SeongUk Park, and Nojun Kwak. Paraphrasing complex network: Network compression via factor transfer. In *Advances in Neural Information Processing Systems*, pp. 2760–2769, 2018. [3](#), [12](#), [14](#)
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009. [5](#)
- Xu Lan, Xiatian Zhu, and Shaogang Gong. Knowledge distillation by on-the-fly native ensemble. In *Advances in neural information processing systems*, 2018. [6](#), [7](#)
- Jinyu Li, Rui Zhao, Jui-Ting Huang, and Yifan Gong. Learning small-size dnn with output-distribution-based criteria. In *Fifteenth annual conference of the international speech communication association*, 2014. [2](#)
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. [1](#), [3](#), [4](#)
- Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3967–3976, 2019. [12](#), [14](#)
- Nikolaos Passalis and Anastasios Tefas. Learning deep representations with probabilistic knowledge transfer. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 268–284, 2018. [12](#), [14](#)
- Baoyun Peng, Xiao Jin, Jiaheng Liu, Shunfeng Zhou, Yichao Wu, Yu Liu, Dongsheng Li, and Zhaoning Zhang. Correlation congruence for knowledge distillation. *arXiv preprint arXiv:1904.01802*, 2019. [12](#), [13](#)
- Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014. [1](#), [3](#)
- Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4510–4520, 2018. [13](#)
- Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *European Conference on Computer Vision*, 2012. [5](#)
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. [13](#)
- Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, Mark Sandler, Andrew Howard, and Quoc V Le. Mnasnet: Platform-aware neural architecture search for mobile. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2820–2828, 2019. [13](#)
- Frederick Tung and Greg Mori. Similarity-preserving knowledge distillation. *arXiv preprint arXiv:1907.09682*, 2019. [12](#), [13](#)
- Junho Yim, Donggyu Joo, Jihoon Bae, and Junmo Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4133–4141, 2017. [3](#), [6](#), [12](#), [14](#)
- Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. *arXiv preprint arXiv:1612.03928*, 2016a. [1](#), [2](#), [3](#), [6](#), [7](#), [12](#), [13](#), [15](#)

Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016b. 5, 13

Richard Zhang, Phillip Isola, and Alexei A Efros. Split-brain autoencoders: Unsupervised learning by cross-channel prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 8

Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6848–6856, 2018. 13

## 6 APPENDIX

### 6.1 OTHER METHODS

We compare to the following other methods from the literature:

1. Knowledge Distillation (KD) (Hinton et al., 2015)
2. Fitnets: Hints for thin deep nets (Ba & Caruana, 2014)
3. Attention Transfer (AT) (Zagoruyko & Komodakis, 2016a)
4. Similarity-Preserving Knowledge Distillation (SP) (Tung & Mori, 2019);
5. Correlation Congruence (CC) (Peng et al., 2019)
6. Variational information distillation for knowledge transfer (VID) (Ahn et al., 2019)
7. Relational Knowledge Distillation (RKD) (Park et al., 2019)
8. Learning deep representations with probabilistic knowledge transfer (PKT) (Passalis & Tefas, 2018)
9. Knowledge transfer via distillation of activation boundaries formed by hidden neurons (AB) (Heo et al., 2019)
10. Paraphrasing complex network: Network compression via factor transfer (FT) (Kim et al., 2018)
11. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning (FSP) (Yim et al., 2017)
12. Like what you like: Knowledge distill via neuron selectivity transfer (NST) (Huang & Wang, 2017)

### 6.2 CONTRASTIVE LOSS – DETAILS

#### 6.2.1 PROOF THAT $h^*(T, S) = q(C = 1|T, S)$

We wish to model some true distribution  $q(C|T = t, S = s)$ .  $C$  is a binary variable, so we can model  $q(C|T = t, S = s)$  as a Bernoulli distribution with a single parameter  $h(S = s, T = t) \in [0, 1]$ , defining for convenience  $h'(C = 1, S = s, T = t) = h(S = s, T = t)$  and  $h'(C = 0, S = s, T = t) = 1 - h(S = s, T = t)$ . The log likelihood function is:

$$\mathbb{E}_{c \sim q(C|S=s, T=t)}[\log h'(C = c, S = s, T = t)] \quad (21)$$

By Gibb’s inequality, the max likelihood fit is  $h'(C = c, S = s, T = t) = q(C = c|S = s, T = t)$ , which also implies that  $h(S = s, T = t) = q(C = 1|S = s, T = t)$ .

We now demonstrate that our objective in Eq. (10) is proportional to a summation over terms Eq. (21) for all  $s \in \mathcal{S}$  and  $t \in \mathcal{T}$ .

$$\mathbb{E}_{s,t \sim q(S,T)} [\mathbb{E}_{c \sim q(C|S=s,T=t)} [\log h'(C = c, S = s, T = t)]] \quad (22)$$

$$= \mathbb{E}_{c,s,t \sim q(C,S,T)} [\log h'(C = c, S = s, T = t)] \quad (23)$$

$$= \mathbb{E}_{s,t \sim q(S,T|C=1)q(C=1)} [\log h(S = s, T = t)] + \mathbb{E}_{s,t \sim q(S,T|C=0)q(C=0)} [\log(1 - h(S = s, T = t))] \quad (24)$$

$$= \frac{1}{N+1} \mathbb{E}_{s,t \sim q(S,T|C=1)} [\log h(S = s, T = t)] + \frac{N}{N+1} \mathbb{E}_{s,t \sim q(S,T|C=0)} [\log(1 - h(S = s, T = t))] \quad (25)$$

Notice that (25) is proportional to Eq. (10) from the main paper. For sufficiently expressive  $h$ , then, each term inside the expectation in Eq. (22) can be maximized, resulting in  $h^*(T = t, S = s) = q(C = 1|T = t, S = s)$  for all  $s$  and  $t$ .  $\square$

### 6.3 NETWORK ARCHITECTURES

*Wide Residual Network (WRN)* (Zagoruyko & Komodakis, 2016b). WRN-d-w represents wide resnet with depth  $d$  and width factor  $w$ .

*resnet* (He et al., 2016). resnet-d represents **cifar**-style resnet with 3 groups of basic blocks with 16, 32, and 64 channels. In our experiments, resnet8 x4 and resnet32 x4 indicate a 4 times wider network (namely, with 64, 128, and 256 channels for each of the residual block)

*ResNet* (He et al., 2016). ResNet-d represents **ImageNet**-style ResNet with Bottleneck blocks and more channels.

*MobileNetV2* Sandler et al. (2018). In our experiments, we use a width multiplier of 0.5.

*vgg* (Simonyan & Zisserman, 2014). the vgg net used in our experiments are adapted from its original ImageNet counterpart.

*ShuffleNetV1* (Zhang et al., 2018), *ShuffleNetV2* (Tan et al., 2019). ShuffleNets are proposed for efficient training and we adapt them to input of size 32x32.

### 6.4 IMPLEMENTATION DETAILS

All methods evaluated in our experiments use SGD.

For CIFAR-100, we initialize the learning rate as 0.05, and decay it by 0.1 every 30 epochs after the first 150 epochs until the last 240 epoch. For MobileNetV2, ShuffleNetV1 and ShuffleNetV2, we use a learning rate of 0.01 as this learning rate is optimal for these models in a grid search, while 0.05 is optimal for other models.

For ImageNet, we follow the standard PyTorch practice but train for 10 more epochs. Batch size is 64 for CIFAR-100 or 256 for ImageNet.

The student is trained by a combination of cross-entropy classification objective and a knowledge distillation objective, shown as follows:

$$\mathcal{L} = \mathcal{L}_{cross-entropy} + \beta \mathcal{L}_{distill} \quad (26)$$

For the weight balance factor  $\beta$ , we directly use the optimal value from the original paper if it is specified, or do a grid search with teacher WRN-40-2 and student WRN-16-2. This results in the following list of  $\beta$  used for different objectives:

1. Fitnets (Ba & Caruana, 2014):  $\beta = 100$
2. AT (Zagoruyko & Komodakis, 2016a):  $\beta = 1000$
3. SP (Tung & Mori, 2019):  $\beta = 3000$
4. CC (Peng et al., 2019):  $\beta = 0.02$

5. VID (Ahn et al., 2019):  $\beta = 1$
6. RKD (Park et al., 2019):  $\beta_1 = 25$  for distance and  $\beta_2 = 50$  for angle. For this loss, we combine both term following the original paper.
7. PKT (Passalis & Tefas, 2018):  $\beta = 30000$
8. AB (Heo et al., 2019):  $\beta = 0$ , distillation happens in a separate pre-training stage where only distillation objective applies.
9. FT (Kim et al., 2018):  $\beta = 500$
10. FSP (Yim et al., 2017):  $\beta = 0$ , distillation happens in a separate pre-training stage where only distillation objective applies.
11. NST (Huang & Wang, 2017):  $\beta = 50$
12. CRD:  $\beta = 0.8$ , in general  $\beta \in [0.5, 1.5]$  works reasonably well.

For KD(Hinton et al., 2015), we follow Eq. 20 and set  $\alpha = 0.9$  and  $T = 4$ .

### 6.5 VISUALIZATION OF THE CORRELATION DISCREPANCY

We visualize the correlation discrepancy for different distillation objectives across various combinations of student and teacher networks. As shown in Fig. 6, our contrastive distillation objective significantly outperforms other objectives, in terms of minimizing the correlation discrepancy between student and teacher networks. The normalized correlation coefficients are computed at the logit layer.

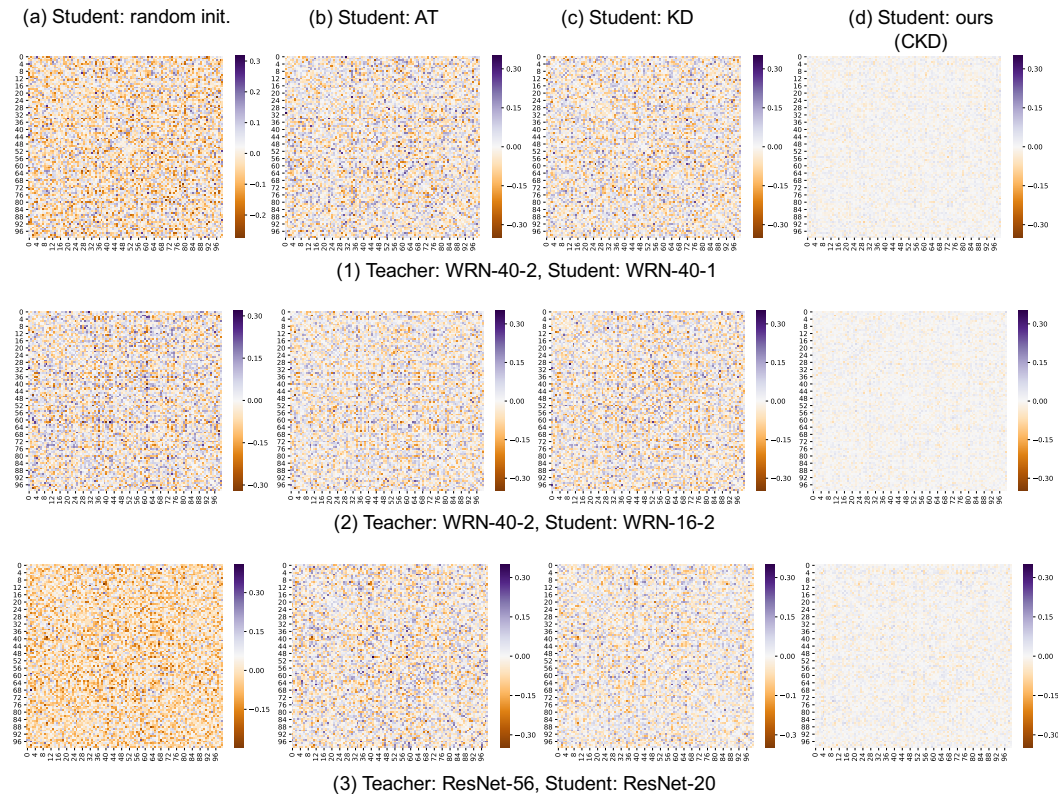


Figure 6: The correlation coefficients between class logits output by the teacher network shows the “dark knowledge” Hinton et al. (2015) that must be transferred to a student networks. A student network that captures these correlations tends to perform better at the task. We visualize here the difference of normalized cross-correlation matrices of the student and teacher at the logits, for different student networks on a Cifar100 knowledge distillation task: (a) A student trained from random initialization; (b) A student distilled by attention transfer Zagoruyko & Komodakis (2016a) (c) A student distilled by KL divergence Hinton et al. (2015); (d) A student distilled by our contrastive objective. Our objective greatly improves the structured knowledge (correlations) in the output units.