

# Hierarchical Self-supervised Augmented Knowledge Distillation

Chuanguang Yang<sup>1,2</sup> Zhulin An<sup>1</sup> Linhang Cai<sup>1,2</sup> Yongjun Xu<sup>1</sup>

<sup>1</sup>Institute of Computing Technology, Chinese Academy of Sciences

<sup>2</sup>University of Chinese Academy of Sciences

IJCAI, 19th-26th August 2021



中科院计算所  
INSTITUTE OF COMPUTING TECHNOLOGY CAS



# Table of Contents

① Introduction

② Methodology

③ Experiments

④ Conclusion



中科院计算所  
INSTITUTE OF COMPUTING TECHNOLOGY, CAS



# Table of Contents

1 Introduction

2 Methodology

3 Experiments

4 Conclusion



中科院计算所  
INSTITUTE OF COMPUTING TECHNOLOGY CAS



The current pattern of KD can be summarized as two critical aspects:

- ④ what kind of knowledge encapsulated in teacher network can be explored for KD:
- ② How to effectively transfer knowledge from teacher to student



中科院计算所  
INSTITUTE OF COMPUTING TECHNOLOGY CAS



## Aspect 1

what kind of knowledge encapsulated in teacher network can be explored for KD:

- ④ **Logit-based class posterior distributions** (Hinton's KD [1])
- ② **Feature-based information:** feature-maps (FitNet [2]), attention maps (AT [3]), gram matrix (FSP [4]), activation boundaries (AB [5]) and so on.
- ③ **Cross-sample relation information:** distance and angle relation (RKD [6]), correlation (CC [7]), similarity (SP [1]), contrastive representations (CRD [2]), and self-supervised contrastive relations (SSKD [3]).



## Previous state-of-the-art: SSKD [3]

Inspired by SimCLR [4], SSKD [3] applies contrastive learning by forcing the image and its transformed version closed against other negative images. It defines the contrastive relationship as knowledge.

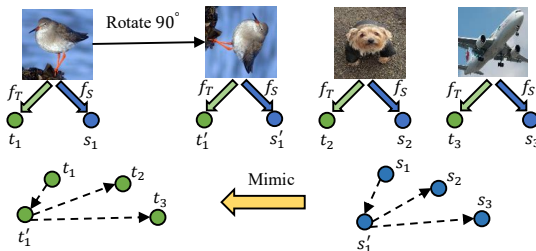


Figure: Self-supervised contrastive relationship [3].



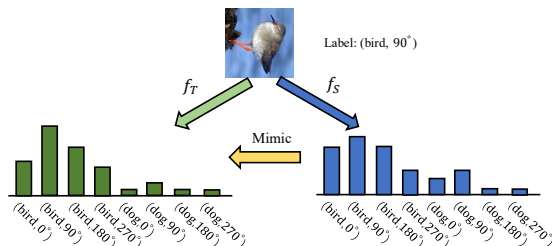
# Introduction

## Motivation

SSKD [3] learns invariant feature representations among transformed images using a self-supervised pretext task with random rotations from  $0^\circ, 90^\circ, 180^\circ, 270^\circ$ , which may destroy the original visual semantics.

## A new type knowledge with self-supervised auxiliary tasks

We introduce a self-supervised augmented distribution that encapsulates the unified knowledge of the original classification task and auxiliary self-supervised task as the richer dark knowledge for the field of KD.



## Empirical verification

We conduct initial exploratory experiments to train a ResNet-18 using rotation as a data augmentation (DA) and a self-supervised augmented label (SAL) as follows. The good performance by SAL further motivates us to define the self-supervised augmented distribution as a promising knowledge for KD.

| Dataset      | Baseline | +DA (Rotation)             | +SAL (Rotation)                  |
|--------------|----------|----------------------------|----------------------------------|
| CIFAR-100    | 78.01    | 77.75( $\downarrow-0.26$ ) | <b>79.76</b> ( $\uparrow+1.75$ ) |
| TinyImageNet | 63.69    | 62.66( $\downarrow-1.03$ ) | <b>65.81</b> ( $\uparrow+2.12$ ) |





## Aspect 2

How to effectively transfer knowledge from teacher to student:

- ① **KD from logits from the final layer:** Hinton's KD [1]
- ② **KD from intermediate feature-maps:** FitNet [2], AT [3], FSP [4] and AB [5].
- ③ **KD from highly abstract feature embeddings before the penultimate layers:** RKD [6], CC [7], SP [1], CRD [2] and SSKD [3].
- ④ **KD assisted with an extra teacher:** HKD [5] introduces an extra teacher model to further bridge the knowledge gap.



中科院计算所  
INSTITUTE OF COMPUTING TECHNOLOGY, CAS



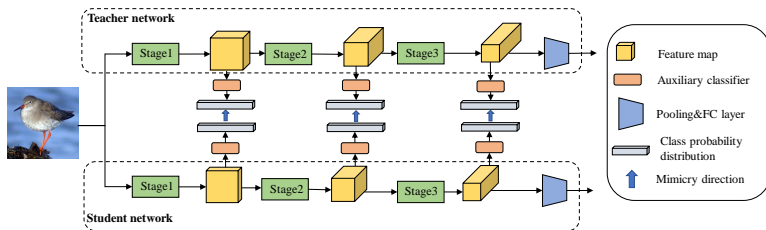
# Introduction

## Motivation

Compared with feature information, the probability distribution is indeed a more robust knowledge for KD [2]. However, it is difficult to explicitly derive probability distributions from hidden layers over the original architecture.

## Our method

Therefore a natural idea is to append several auxiliary classifiers to the network at various hidden layers to generate multi-level probability distributions from hierarchical feature maps.



# Table of Contents

1 Introduction

2 Methodology

3 Experiments

4 Conclusion



中科院计算所  
INSTITUTE OF COMPUTING TECHNOLOGY CAS



**CNN Architecture:** A CNN can be decomposed into a feature extractor  $\Phi(\cdot; \boldsymbol{\mu})$  and a linear classifier  $g(\cdot; \boldsymbol{w})$ , where  $\boldsymbol{\mu}$  and  $\boldsymbol{w}$  are weight tensors.

## Conventional Class Probability Distribution

Given an input sample  $\mathbf{x} \in \mathcal{X}$ ,  $\mathcal{X}$  is the training set,  $\mathbf{z} = \Phi(\mathbf{x}; \boldsymbol{\mu}) \in \mathbb{R}^d$  is the extracted feature embedding vector, where  $d$  is the embedding size. We consider a conventional  $N$ -way object classification task with the label space  $\mathcal{N} = \{1, \dots, N\}$ . The linear classifier maps the feature embedding  $\mathbf{z}$  to a predictive class probability distribution  $\mathbf{p}(\mathbf{x}; \tau) = \sigma(g(\mathbf{z}; \boldsymbol{w})/\tau) \in \mathbb{R}^N$ .

Hinton's KD [1] uses conventional class probability distribution as knowledge.



中科院计算所  
INSTITUTE OF COMPUTING TECHNOLOGY, CAS



## Self-supervised Augmented Label Space

Assuming that we define  $M$  various image transformations  $\{t_j\}_{j=1}^M$  with the label space  $\mathcal{M} = \{1, \dots, M\}$ , where  $t_1(\mathbf{x}) = \mathbf{x}$ . The label space of this task is  $\mathcal{K} = \mathcal{N} \otimes \mathcal{M}$ , here  $\otimes$  is the Cartesian product.  $|\mathcal{K}| = N * M$ , where  $|\cdot|$  is the cardinality of the label collection,  $*$  denotes element-wise multiplication.

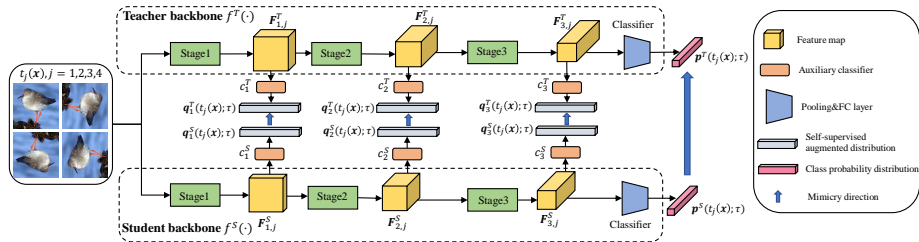
## Self-supervised Augmented Distribution

Given a transformed sample  $\tilde{\mathbf{x}} \in \{t_j(\mathbf{x})\}_{j=1}^M$  by applying one transformation on  $\mathbf{x}$ ,  $\tilde{\mathbf{z}} = \Phi(\tilde{\mathbf{x}}; \boldsymbol{\mu}) \in \mathbb{R}^d$  is the extracted feature embedding vector,  $\mathbf{q}(\tilde{\mathbf{x}}; \tau) = \sigma(g(\tilde{\mathbf{z}}; \mathbf{w})/\tau) \in \mathbb{R}^{N*M}$  is the predictive distribution over the joint label space  $\mathcal{K}$ , where weight tensor  $\mathbf{w} \in \mathbb{R}^{(N*M) \times d}$ .



## Hierarchical Self-supervised Augmented Knowledge Distillation

We guide all auxiliary classifiers attached to the original network to learn informative self-supervised augmented distributions. Furthermore, we perform knowledge distillation between teacher and student towards all auxiliary classifiers in a one-to-one manner.



## Pre-train a teacher network

We denote the teacher network as  $f^T(\cdot)$  and  $L$  auxiliary classifiers as  $\{c_l^T(\cdot)\}_{l=1}^L$ .

- ① we train the  $f^T(\cdot)$  with normal data  $\mathbf{x}$  by the conventional Cross-Entropy (CE) loss to fit the ground-truth label  $y \in \mathcal{N}$ .
- ② we aim to train  $L$  auxiliary classifiers  $\{c_l^T(\cdot)\}_{l=1}^L$  for learning self-supervised augmented labels  $k_j$ .

The overall loss for training a teacher is shown in Eq. (1).

$$\mathcal{L}_T = \mathbb{E}_{x \in \mathcal{X}} [\mathcal{L}_{ce}(p^T(x; \tau), y) + \frac{1}{M} \sum_{j=1}^M \sum_{l=1}^L \mathcal{L}_{ce}(q_l^T(t_j(x); \tau), k_j)] \quad (1)$$

$p^T(x; \tau) = \sigma(f^T(x)/\tau) \in \mathbb{R}^N$  is predictive class probability distribution,  
 $q_l^T(t_j(x); \tau) = \sigma(c_l^T(F_{l,j}^T))/\tau \in \mathbb{R}^{N \times M}$  is self-supervised augmented distributions .

## Train a student network supervised by a teacher network

We denote the student backbone network as  $f^S(\cdot)$  and  $L$  auxiliary classifiers as  $\{c_l^S(\cdot)\}_{l=1}^L$ . The overall loss includes a task loss from ground-truth labels and mimicry losses from the pre-trained teacher.

- ④ Task loss from ground-truth labels  $\mathcal{L}_{task}$
- ② Mimicry loss from self-supervised augmented distributions  $\mathcal{L}_{kl-q}$
- ③ Mimicry loss from class probability distributions  $\mathcal{L}_{kl-p}$





## Task loss from ground-truth labels $\mathcal{L}_{task}$

We force the  $f^S(\cdot)$  to fit the normal data  $x$  with the ground-truth  $y$  as the task loss:

$$\mathcal{L}_{task} = \mathcal{L}_{ce}(p^S(x; \tau), y) \quad (2)$$

Where  $p^S(x; \tau) = \sigma(f^S(x)/\tau) \in \mathbb{R}^N$  is the predictive class probability distribution.



## Mimicry loss from self-supervised augmented distributions $\mathcal{L}_{kl-q}$

We consider transferring hierarchical self-supervised augmented distributions  $\{q_l^T(t_j(x); \tau)\}_{l=1}^L$  generated from  $L$  auxiliary classifiers of the teacher network to corresponding  $\{q_l^S(t_j(x); \tau)\}_{l=1}^L$  generated from  $L$  auxiliary classifiers of the student network, respectively. The transfer performs in a one-to-one manner by KL-divergence loss  $D_{KL}$ .

$$\mathcal{L}_{kl-q} = \frac{1}{M} \sum_{j=1}^M \sum_{l=1}^L \tau^2 D_{KL}(q_l^T(t_j(x); \tau) \parallel q_l^S(t_j(x); \tau)) \quad (3)$$



## Mimicry loss from class probability distributions $\mathcal{L}_{kl-p}$

We transfer the original class probability distributions generated from the final layer between teacher and student. Specifically, we transfer the knowledge derived from both the normal and transformed data  $\{t_j(x)\}_{j=1}^M$ , where  $t_1(x) = x$ .

$$\mathcal{L}_{kl-p} = \frac{1}{M} \sum_{j=1}^M \tau^2 D_{\text{KL}}(p^T(t_j(x); \tau) \parallel p^S(t_j(x); \tau)) \quad (4)$$



## Overall loss for training the student network

We summarize the task loss and mimicry loss as the overall loss  $\mathcal{L}_S$  for training the student network:

$$\mathcal{L}_S = \mathbb{E}_{x \in \mathcal{X}} [\mathcal{L}_{task} + \mathcal{L}_{kl\_q} + \mathcal{L}_{kl\_p}] \quad (5)$$

Following the wide practice, we set the hyper-parameter  $\tau = 1$  in task loss and  $\tau = 3$  in mimicry loss. Besides, we do not introduce other hyper-parameters.



# Table of Contents

1 Introduction

2 Methodology

3 Experiments

4 Conclusion



中科院计算所  
INSTITUTE OF COMPUTING TECHNOLOGY CAS



# Experiments

## Experiments on CIFAR-100

Our HSAKD significantly outperforms the best-competing method SSKD across all network pairs with an average accuracy gain of 2.56%

| Teacher  | WRN-40-2     | WRN-40-2     | ResNet56     | ResNet32×4   | VGG13        | ResNet50     | WRN-40-2     | ResNet32×4   |
|----------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Student  | WRN-16-2     | WRN-40-1     | ResNet20     | ResNet8×4    | MobileNetV2  | MobileNetV2  | ShuffleNetV1 | ShuffleNetV2 |
| Teacher  | 76.45        | 76.45        | 73.44        | 79.63        | 74.64        | 79.34        | 76.45        | 79.63        |
| Teacher* | 80.70        | 80.70        | 77.20        | 83.73        | 78.48        | 83.85        | 80.70        | 83.73        |
| Student  | 73.57(±0.23) | 71.95(±0.59) | 69.62(±0.26) | 72.95(±0.24) | 73.51(±0.26) | 73.51(±0.26) | 71.74(±0.35) | 72.96(±0.33) |
| KD       | 75.23(±0.23) | 73.90(±0.44) | 70.91(±0.10) | 73.54(±0.26) | 75.21(±0.24) | 75.80(±0.46) | 75.83(±0.18) | 75.43(±0.33) |
| FitNet   | 75.30(±0.42) | 74.30(±0.42) | 71.21(±0.16) | 75.37(±0.12) | 75.42(±0.34) | 75.41(±0.07) | 76.27(±0.18) | 76.91(±0.06) |
| AT       | 75.64(±0.31) | 74.32(±0.23) | 71.35(±0.09) | 75.06(±0.19) | 74.08(±0.21) | 76.57(±0.20) | 76.51(±0.44) | 76.32(±0.12) |
| AB       | 71.26(±1.32) | 74.55(±0.46) | 71.56(±0.19) | 74.31(±0.09) | 74.98(±0.44) | 75.87(±0.39) | 76.43(±0.09) | 76.40(±0.29) |
| VID      | 75.31(±0.22) | 74.23(±0.28) | 71.35(±0.09) | 75.07(±0.35) | 75.67(±0.13) | 75.97(±0.08) | 76.24(±0.44) | 75.98(±0.41) |
| RKD      | 75.33(±0.14) | 73.90(±0.26) | 71.67(±0.08) | 74.17(±0.22) | 75.54(±0.36) | 76.20(±0.06) | 75.74(±0.32) | 75.42(±0.25) |
| SP       | 74.35(±0.59) | 72.91(±0.24) | 71.45(±0.38) | 75.44(±0.11) | 75.68(±0.35) | 76.35(±0.14) | 76.40(±0.37) | 76.43(±0.21) |
| CC       | 75.30(±0.03) | 74.46(±0.05) | 71.44(±0.10) | 74.40(±0.24) | 75.66(±0.33) | 76.05(±0.25) | 75.63(±0.30) | 75.74(±0.18) |
| CRD      | 75.81(±0.33) | 74.76(±0.25) | 71.83(±0.42) | 75.77(±0.24) | 76.13(±0.16) | 76.89(±0.27) | 76.37(±0.23) | 76.51(±0.09) |
| SSKD     | 76.16(±0.17) | 75.84(±0.04) | 70.80(±0.02) | 75.83(±0.29) | 76.21(±0.16) | 78.21(±0.16) | 76.71(±0.31) | 77.64(±0.24) |
| Ours     | 77.20(±0.17) | 77.00(±0.21) | 72.58(±0.33) | 77.26(±0.14) | 77.45(±0.21) | 78.79(±0.11) | 78.51(±0.20) | 79.93(±0.11) |
| Ours*    | 78.67(±0.20) | 78.12(±0.25) | 73.73(±0.10) | 77.69(±0.05) | 79.27(±0.12) | 79.43(±0.24) | 80.11(±0.32) | 80.86(±0.15) |

Table: Top-1 accuracy (%) comparison of SOTA methods.



中科院计算所  
INSTITUTE OF COMPUTING TECHNOLOGY CAS



## Experiments on ImageNet

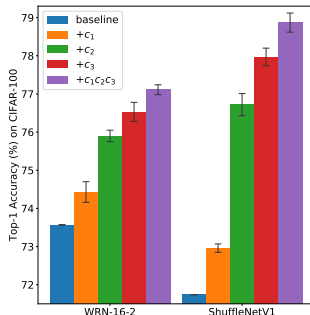
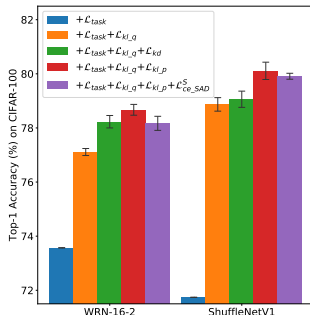
Our HSAKD significantly outperforms the best-competing method SSKD on the pair of ResNet-34 and ResNet-18 with a top-1 accuracy gain of 0.77%

| Teacher   | Student   | Acc   | Teacher | Teacher* | Student | KD    | AT    | CC    | SP    | RKD   | CRD   | SSKD         | Ours  | Ours*        |
|-----------|-----------|-------|---------|----------|---------|-------|-------|-------|-------|-------|-------|--------------|-------|--------------|
| ResNet-34 | ResNet-18 | Top-1 | 73.31   | 75.48    | 69.75   | 70.66 | 70.70 | 69.96 | 70.62 | 71.34 | 71.38 | <u>71.62</u> | 72.16 | <b>72.39</b> |
|           |           | Top-5 | 91.42   | 92.67    | 89.07   | 89.88 | 90.00 | 89.17 | 89.80 | 90.37 | 90.49 | <u>90.67</u> | 90.85 | <b>91.00</b> |

Table: Top-1 accuracy (%) comparison on ImageNet.

## Ablation study on CIFAR-100

- **Effect of loss terms (*left*):** each loss term is indispensable, and  $\mathcal{L}_{kl\_q}$  contributes most by distilling the self-supervised augmented distributions.
- **Effect of auxiliary classifiers (*right*):** each auxiliary classifier is indispensable. The auxiliary classifier attached in the deeper layer often achieves more accuracy gains than that in the shallower layer. Using all auxiliary classifiers can maximize accuracy gains.



中国科学院  
计算技术研究所





## Transfer Experiments on STL-10 and Tiny ImageNet

Our HSAKD can significantly outperform the best-competing SSKD by 3.63% on STL-10 and 3.50% on TinyImageNet.

| Transferred Dataset     | Baseline | KD    | FitNet | AT    | AB    | VID   | RKD   | SP    | CC    | CRD   | SSKD         | Ours         |
|-------------------------|----------|-------|--------|-------|-------|-------|-------|-------|-------|-------|--------------|--------------|
| CIFAR-100→ STL-10       | 67.76    | 67.90 | 69.41  | 67.37 | 67.82 | 69.29 | 69.74 | 68.96 | 69.13 | 70.09 | <u>71.03</u> | <b>74.66</b> |
| CIFAR-100→ TinyImageNet | 34.69    | 34.15 | 36.04  | 34.44 | 34.79 | 36.09 | 37.21 | 35.69 | 36.43 | 38.17 | <u>39.07</u> | <b>42.57</b> |

**Table:** Linear classification accuracy (%) of transfer learning on the student MobileNetV2 pre-trained using the teacher VGG-13.

## Transfer Experiments on Pascal VOC

Our HSAKD outperforms the original baseline by 2.27% mAP and the best-competing SSKD by 0.85% mAP on downstream object detection.

| Baseline | KD    | CRD   | SSKD         | Ours         |
|----------|-------|-------|--------------|--------------|
| 76.18    | 77.06 | 77.36 | <u>77.60</u> | <b>78.45</b> |

**Table:** Comparison of detection mAP (%) on Pascal VOC using ResNet-18 as the backbone pre-trained by various KD methods.



## Efficacy under Few-shot Scenario

Our HSAKD can consistently surpass other KD methods by large margins under various few-shot settings.

| Percentage | KD                  | CRD                         | SSKD                        | Ours                        |
|------------|---------------------|-----------------------------|-----------------------------|-----------------------------|
| 25%        | 65.15( $\pm 0.23$ ) | 65.80( $\pm 0.61$ )         | <u>67.82</u> ( $\pm 0.30$ ) | <b>68.50</b> ( $\pm 0.24$ ) |
| 50%        | 68.61( $\pm 0.22$ ) | 69.91( $\pm 0.20$ )         | <u>70.08</u> ( $\pm 0.13$ ) | <b>72.18</b> ( $\pm 0.41$ ) |
| 75%        | 70.34( $\pm 0.09$ ) | <u>70.98</u> ( $\pm 0.43$ ) | 70.47( $\pm 0.14$ )         | <b>73.26</b> ( $\pm 0.11$ ) |

**Table:** Top-1 accuracy (%) comparison on CIFAR-100 under few-shot scenario with various percentages of training samples. We use the ResNet56-ResNet20 as the teacher-student pair for evaluation.

# Table of Contents

1 Introduction

2 Methodology

3 Experiments

4 Conclusion



中科院计算所  
INSTITUTE OF COMPUTING TECHNOLOGY, CAS



- ④ **Knowledge definition:** We introduce a self-supervised augmented distribution that encapsulates the unified knowledge of the original classification task and auxiliary self-supervised task as the richer dark knowledge for the field of KD.
- ② **Knowledge transfer:** We propose a one-to-one probabilistic knowledge distillation framework by leveraging the architectural auxiliary classifiers, facilitating comprehensive knowledge transfer and alleviating the mismatch problem of abstraction levels when existing a large architecture gap.
- ③ **Experimental results:** HSAKD significantly refreshes the results achieved by previous SOTA SSKD on standard image classification benchmarks. It can also learn well-general feature representations for downstream semantic recognition tasks.



# Reference



Geoffrey Hinton, Oriol Vinyals, and Jeff Dean.

Distilling the knowledge in a neural network.  
*arXiv preprint arXiv:1503.02531*, 2015.



Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio.

Fitnets: Hints for thin deep nets.  
*ICLR*, 2015.



Sergey Zagoruyko and Nikos Komodakis.

Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer.  
*ICLR*, 2017.



Junho Yim, Donggyu Joo, Ji-Hoon Bae, and Junmo Kim.

A gift from knowledge distillation: Fast optimization, network minimization and transfer learning.  
In *CVPR*, pages 7130–7138, 2017.



Byeongho Heo, Minsik Lee, Sangdoo Yun, and Jin Young Choi.

Knowledge transfer via distillation of activation boundaries formed by hidden neurons.  
In *AAAI*, volume 33, pages 3779–3787, 2019.



Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho.

Relational knowledge distillation.  
In *CVPR*, pages 3967–3976, 2019.



Baoyun Peng, Xiao Jin, Jiaheng Liu, Dongsheng Li, Yichao Wu, Yu Liu, Shunfeng Zhou, and Zhaoning Zhang.

Correlation congruence for knowledge distillation.  
In *ICCV*, pages 5007–5016, 2019.



中科院计算所  
INSTITUTE OF COMPUTING TECHNOLOGY CAS





Frederick Tung and Greg Mori.  
Similarity-preserving knowledge distillation.  
In *ICCV*, pages 1365–1374, 2019.



Yonglong Tian, Dilip Krishnan, and Phillip Isola.  
Contrastive representation distillation.  
*ICLR*, 2020.



Guodong Xu, Ziwei Liu, Xiaoxiao Li, and Chen Change Loy.  
Knowledge distillation meets self-supervision.  
In *ECCV*, pages 588–604, 2020.



Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton.  
A simple framework for contrastive learning of visual representations.  
*ICML*, pages 1597–1607, 2020.



Nikolaos Passalis, Maria Tzelepi, and Anastasios Tefas.  
Heterogeneous knowledge distillation using information flow modeling.  
In *CVPR*, pages 2339–2348, 2020.



Gidaris, Spyros and Singh, Praveer and Komodakis, Nikos.  
Unsupervised representation learning by predicting image rotations.  
In *ICLR*, 2018

