

Relational Representation Distillation

Supplementary Material

1. Algorithm

Algorithm 1 provides the pseudo-code of RRD.

2. Implementation Details

We implement RRD in PyTorch following the implementation of CRD¹.

2.1. Baseline Methods

We compare our approach to the following state-of-the-art methods from the literature: (1) Knowledge Distillation (KD) [11]; (2) FitNets: Hints for Thin Deep Nets [22]; (3) Attention Transfer (AT) [29]; (4) Similarity-Preserving Knowledge Distillation (SP) [26]; (5) Correlation Congruence (CC) [21]; (6) Variational Information Distillation for Knowledge Transfer (VID) [1]; (7) Relational Knowledge Distillation (RKD) [19]; (8) Learning Deep Representations with Probabilistic Knowledge Transfer (PKT) [20]; (9) Knowledge Transfer via Distillation of Activation Boundaries Formed by Hidden Neurons (AB) [10]; (10) Paraphrasing Complex Network: Network Compression via Factor Transfer (FT) [13]; (11) A Gift from Knowledge Distillation: Fast Optimization, Network Minimization and Transfer Learning (FSP) [28]; (12) Like What You Like: Knowledge Distill via Neuron Selectivity Transfer (NST) [12]; (13) Contrastive Representation Distillation (CRD) [25]; (14) A Comprehensive Overhaul of Feature Distillation (OFD); (15) Rethinking Soft Labels for Knowledge Distillation: A Bias-Variance Tradeoff Perspective (WSLD) [32]; (16) Respecting Transfer Gap in Knowledge Distillation (IPWD) [18]; (17) Knowledge Distillation via Softmax Regression Representation Learning (SRRL) [27]; (18) Cross-Layer Distillation with Semantic Calibration (SemCKD) [2]; (19) Distilling Knowledge via Knowledge Review (ReviewKD) [5]; (20) Knowledge Distillation with the Reused Teacher Classifier (SimKD) [3]; (21) Searching A Fast Knowledge Distillation Process via Meta Optimization (DistPro) [6]; (22) Knowledge Distillation via N-to-One Representation Matching (NORM) [14]; (23) Information Theoretic Representation (ITRD) [17]; (24) Feature Kernel Distillation (FKD) [7]; (25) Complementary Relation Contrastive Distillation (CRCDD) [33]; (26) Distilling Knowledge from Self-Supervised Teacher by Embedding Graph Alignment (EGA) [16]; (27) Wasserstein Contrastive Representation Distillation (WCoRD) [4].

¹Available at: <https://github.com/HobbitLong/RepDistiller>.

2.2. Network Architectures

We use the following network architectures as described in [25]: (1) Wide Residual Network (WRN) [30], where WRN- $d-w$ represents a wide ResNet with depth d and width factor w ; (2) ResNet [8], where resnet- d represents a CIFAR-style ResNet with 3 groups of basic blocks having 16, 32, and 64 channels, respectively, and resnet- 8×4 and resnet- 32×4 indicate a 4-times wider network with 64, 128, and 256 channels; (3) ResNet [8], where ResNet- d represents an ImageNet-style ResNet with Bottleneck blocks and more channels; (4) MobileNet-v2 [23], using a width multiplier of 0.5 in our experiments; (5) VGG [24], where the VGG network used is adapted from its original ImageNet counterpart; and (6) ShuffleNet-v1 [31] and ShuffleNet-v2 [15], which are adapted for efficient training with input sizes of 32×32 .

2.3. Optimization

All methods evaluated in our experiments use SGD with 0.9 Nesterov momentum. 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 MobileNet-v2, ShuffleNet-v1, and ShuffleNet-v2, 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. The batch size is set to 64 for CIFAR-100, and the weight decay is set to 5×10^{-4} . For ImageNet, the initial learning rate is set to 0.1 and then divided by 10 at the 30th, 60th, and 90th epochs of the total 120 training epochs. The mini-batch size is set to 256, and the weight decay is set to 1×10^{-4} . All results are reported as means over three trials, except for the results on ImageNet, which are reported in a single trial.

3. Results

Table 1 and Table 2 provide a comprehensive overview of the top-1 accuracies of student networks trained with various state-of-the-art distillation techniques across a wide range of teacher-student architectural combinations. Unlike the main text, which summarizes a subset of results, these tables offer an extended comparison involving more models and training configurations. Our proposed method, RRD, shows strong performance across diverse network architectures and teacher-student pairs. RRD performs nearly as well as the top methods in knowledge distillation, achieving accuracy rates very close to the best-performing techniques, indicating an effective balance between simplicity and performance.

Algorithm 1 Pseudocode of RRD in a PyTorch-like style.

```
# f_t, f_s: outputs at the penultimate layer of teacher and student networks
# t_dim: The input feature dimension for the teacher
# s_dim: The input feature dimension for the student
# feat_dim: The projection feature space dimension
# nce_k: number of instances in queue
# nce_t_s, nce_t_t: the temperature parameters for student and teacher networks
# N: batch size

class RRDLoss(nn.Module):
    def __init__(self, s_dim, t_dim, feat_dim, nce_k=16384, nce_t_t=0.07, nce_t_s=0.04):
        super(RRDLoss, self).__init__()

        # embedding layer
        self.embed_s = nn.Linear(s_dim, feat_dim)
        self.embed_t = nn.Linear(t_dim, feat_dim)

        # memory buffer
        self.register_buffer("queue", torch.randn(nce_k, feat_dim))
        self.queue = F.normalize(self.queue, dim=0)
        self.register_buffer("queue_ptr", torch.zeros(1, dtype=torch.long))

    def forward(self, f_s, f_t):
        f_s = self.embed_s(f_s)
        f_t = self.embed_t(f_t)

        f_s = F.normalize(f_s, dim=1)
        f_t = F.normalize(f_t, dim=1)

        l_s = torch.einsum("nc,kc->nk", [f_s, self.queue])
        l_t = torch.einsum("nc,kc->nk", [f_t, self.queue])

        loss = -torch.sum(
            F.softmax(l_t / self.nce_t_t, dim=1) *
            F.log_softmax(l_s / self.nce_t_s, dim=1), dim=1).mean()

        self._dequeue_and_enqueue(f_t)

        return loss
```

References

- [1] 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*, pages 9163–9171, 2019. 1, 3, 4
- [2] Defang Chen, Jian-Ping Mei, Yuan Zhang, Can Wang, Yan Feng, and Chun Chen. Cross-layer distillation with semantic calibration, 2021. 1, 3, 4
- [3] Defang Chen, Jian-Ping Mei, Hailin Zhang, Can Wang, Yan Feng, and Chun Chen. Knowledge distillation with the reused teacher classifier, 2022. 1, 3, 4
- [4] Liqun Chen, Dong Wang, Zhe Gan, Jingjing Liu, Ricardo Henao, and Lawrence Carin. Wasserstein contrastive representation distillation, 2021. 1, 3, 4
- [5] Pengguang Chen, Shu Liu, Hengshuang Zhao, and Jiaya Jia. Distilling knowledge via knowledge review, 2021. 1, 3, 4
- [6] Xueqing Deng, Dawei Sun, Shawn Newsam, and Peng Wang. Distpro: Searching a fast knowledge distillation process via meta optimization, 2022. 1, 3, 4
- [7] Bobby He and Mete Ozay. Feature kernel distillation. In *International Conference on Learning Representations*, 2022. 1, 3, 4
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015. 1
- [9] Byeongho Heo, Jeesoo Kim, Sangdoon Yun, Hyojin Park, Nojun Kwak, and Jin Young Choi. A comprehensive overhaul of feature distillation, 2019. 3, 4
- [10] Byeongho Heo, Minsik Lee, Seong Joon Yun, Jin Young Choi, and In So Kweon. Knowledge transfer via distillation of activation boundaries formed by hidden neurons. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 3779–3787, 2019. 1, 3, 4
- [11] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network, 2015. 1, 3, 4
- [12] Zehao Huang and Naiyan Wang. Like what you like: Knowledge distill via neuron selectivity transfer. In *Advances in Neural Information Processing Systems*, pages 185–195, 2017. 1, 3, 4
- [13] Jangho Kim, Seongwon Park, and Nojun Kwak. Paraphrasing complex network: Network compression via factor transfer. In *Advances in Neural Information Processing Systems*, pages 2760–2769, 2018. 1, 3, 4
- [14] Xiaolong Liu, Lujun Li, Chao Li, and Anbang Yao. Norm: Knowledge distillation via n-to-one representation matching, 2023. 1, 3, 4
- [15] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 116–131, 2018. 1
- [16] Yuchen Ma, Yanbei Chen, and Zeynep Akata. Distilling knowledge from self-supervised teacher by embedding graph alignment, 2022. 1
- [17] Roy Miles, Adrian Lopez Rodriguez, and Krystian Mikolajczyk. Information theoretic representation distillation, 2022. 1, 3, 4

Table 1. Test top-1 accuracy (%) of student networks on CIFAR-100, comparing students and teachers of the same architecture using various distillation methods. The values in bold indicate the maximum of each column. \uparrow denotes outperformance over KD and \downarrow denotes underperformance.

Teacher	WRN-40-2	WRN-40-2	resnet-56	resnet-110	resnet-110	resnet-32x4	VGG-13
Student	WRN-16-2	WRN-40-1	resnet-20	resnet-20	resnet-32	resnet-8x4	VGG-8
<i>Teacher</i>	75.61	75.61	72.34	74.31	74.31	79.42	74.64
<i>Student</i>	73.26	71.98	69.06	69.06	71.14	72.50	70.36
KD [11]	74.92	73.54	70.66	70.67	73.08	73.33	72.98
FitNet [22]	73.58 (\downarrow)	72.24 (\downarrow)	69.21 (\downarrow)	68.99 (\downarrow)	71.06 (\downarrow)	73.50 (\uparrow)	71.02 (\downarrow)
AT [29]	74.08 (\downarrow)	72.77 (\downarrow)	70.55 (\downarrow)	70.22 (\downarrow)	72.31 (\downarrow)	73.44 (\uparrow)	71.43 (\downarrow)
SP [26]	73.83 (\downarrow)	72.43 (\downarrow)	69.67 (\downarrow)	70.04 (\downarrow)	72.69 (\downarrow)	72.94 (\downarrow)	72.68 (\downarrow)
CC [21]	73.56 (\downarrow)	72.21 (\downarrow)	69.63 (\downarrow)	69.48 (\downarrow)	71.48 (\downarrow)	72.97 (\downarrow)	70.81 (\downarrow)
VID [1]	74.11 (\downarrow)	73.30 (\downarrow)	70.38 (\downarrow)	70.16 (\downarrow)	72.61 (\downarrow)	73.09 (\downarrow)	71.23 (\downarrow)
RKD [19]	73.35 (\downarrow)	72.22 (\downarrow)	69.61 (\downarrow)	69.25 (\downarrow)	71.82 (\downarrow)	71.90 (\downarrow)	71.48 (\downarrow)
PKT [20]	74.54 (\downarrow)	73.45 (\downarrow)	70.34 (\downarrow)	70.25 (\downarrow)	72.61 (\downarrow)	73.64 (\uparrow)	72.88 (\downarrow)
AB [10]	72.50 (\downarrow)	72.38 (\downarrow)	69.47 (\downarrow)	69.53 (\downarrow)	70.98 (\downarrow)	73.17 (\downarrow)	70.94 (\downarrow)
FT [13]	73.25 (\downarrow)	71.59 (\downarrow)	69.84 (\downarrow)	70.22 (\downarrow)	72.37 (\downarrow)	72.86 (\downarrow)	70.58 (\downarrow)
FSP [28]	72.91 (\downarrow)	n/a	69.95 (\downarrow)	70.11 (\downarrow)	71.89 (\downarrow)	72.62 (\downarrow)	70.33 (\downarrow)
NST [12]	73.68 (\downarrow)	72.24 (\downarrow)	69.60 (\downarrow)	69.53 (\downarrow)	71.96 (\downarrow)	73.30 (\downarrow)	71.53 (\downarrow)
CRD [25]	75.48 (\uparrow)	74.14 (\uparrow)	71.16 (\uparrow)	71.46 (\uparrow)	73.48 (\uparrow)	75.51 (\uparrow)	73.94 (\uparrow)
CRD+KD [25]	75.64 (\uparrow)	74.38 (\uparrow)	71.63 (\uparrow)	71.56 (\uparrow)	73.75 (\uparrow)	75.46 (\uparrow)	74.29 (\uparrow)
OFD [9]	75.24 (\uparrow)	74.33 (\uparrow)	70.38 (\downarrow)	n/a	73.23 (\uparrow)	74.95 (\uparrow)	73.95 (\uparrow)
WSLD [32]	n/a	73.74 (\uparrow)	71.53 (\uparrow)	n/a	73.36 (\uparrow)	74.79 (\uparrow)	n/a
IPWD [18]	n/a	74.64 (\uparrow)	71.32 (\uparrow)	n/a	73.91 (\uparrow)	76.03 (\uparrow)	n/a
SRRL [27]	n/a	74.64 (\uparrow)	n/a	n/a	n/a	75.39 (\uparrow)	n/a
SemCKD [2]	n/a	74.41 (\uparrow)	n/a	n/a	n/a	76.23 (\uparrow)	n/a
ReviewKD [5]	76.12 (\uparrow)	75.09 (\uparrow)	71.89 (\uparrow)	n/a	73.89 (\uparrow)	75.63 (\uparrow)	74.84 (\uparrow)
SimKD [3]	n/a	75.56 (\uparrow)	n/a	n/a	n/a	78.08 (\uparrow)	n/a
DistPro [6]	76.36 (\uparrow)	n/a	72.03 (\uparrow)	n/a	73.74 (\uparrow)	n/a	n/a
NORM [14]	75.65 (\uparrow)	74.82 (\uparrow)	71.35 (\uparrow)	71.55 (\uparrow)	73.67 (\uparrow)	76.49 (\uparrow)	73.95 (\uparrow)
NORM+KD [14]	76.26 (\uparrow)	75.42 (\uparrow)	71.61 (\uparrow)	72.00 (\uparrow)	74.95 (\uparrow)	76.98 (\uparrow)	74.46 (\uparrow)
NORM+CRD [14]	76.02 (\uparrow)	75.37 (\uparrow)	71.51 (\uparrow)	71.90 (\uparrow)	73.81 (\uparrow)	76.49 (\uparrow)	73.58 (\uparrow)
WCoRD [4]	75.88 (\uparrow)	74.73 (\uparrow)	71.56 (\uparrow)	71.57 (\uparrow)	73.81 (\uparrow)	75.95 (\uparrow)	74.55 (\uparrow)
WCoRD+KD [4]	76.11 (\uparrow)	74.72 (\uparrow)	71.92 (\uparrow)	71.88 (\uparrow)	74.20 (\uparrow)	76.15 (\uparrow)	74.72 (\uparrow)
CRCD [33]	76.67 (\uparrow)	75.95 (\uparrow)	73.21 (\uparrow)	72.33 (\uparrow)	74.98 (\uparrow)	76.42 (\uparrow)	74.97 (\uparrow)
FKD [7]	n/a	n/a	n/a	n/a	n/a	75.57 (\uparrow)	73.78 (\uparrow)
ITRD (corr) [17]	75.85 (\uparrow)	74.90 (\uparrow)	71.45 (\uparrow)	71.77 (\uparrow)	74.02 (\uparrow)	75.63 (\uparrow)	74.70 (\uparrow)
ITRD (corr+mi) [17]	76.12 (\uparrow)	75.18 (\uparrow)	71.47 (\uparrow)	71.99 (\uparrow)	74.26 (\uparrow)	76.19 (\uparrow)	74.93 (\uparrow)
RRD (ours)	75.01 (\uparrow)	73.55 (\uparrow)	70.71 (\uparrow)	70.72 (\uparrow)	73.10 (\uparrow)	74.48 (\uparrow)	73.99 (\uparrow)
RRD+KD (ours)	75.66 (\uparrow)	74.39 (\uparrow)	72.19 (\uparrow)	71.74 (\uparrow)	73.54 (\uparrow)	75.08 (\uparrow)	74.32 (\uparrow)

[18] Yulei Niu, Long Chen, Chang Zhou, and Hanwang Zhang. Respecting transfer gap in knowledge distillation, 2022. **1, 3, 4**

[19] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3967–3976, 2019. **1, 3, 4**

[20] Nikolaos Passalis and Anastasios Tefas. Learning deep representations with probabilistic knowledge transfer. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 268–284, 2018. **1, 3, 4**

[21] Baoyun Peng, Xi Li, Yifan Wu, Yizhou Fan, Bo Wang, Qi

Tian, and Jun Liang. Correlation congruence for knowledge distillation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5007–5016, 2019. **1, 3, 4**

[22] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. In *Proceedings of the 4th International Conference on Learning Representations*, 2014. **1, 3, 4**

[23] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE*

Table 2. Test top-1 accuracy (%) of student networks on CIFAR-100 involving students and teachers from different architectures, using various distillation methods. The values in bold indicate the maximum of each column. \uparrow denotes outperformance over KD and \downarrow denotes underperformance.

Teacher	VGG-13	ResNet-50	ResNet-50	ResNet-32x4	ResNet-32x4	WRN-40-2
Student	MobileNet-v2	MobileNet-v2	VGG-8	ShuffleNet-v1	ShuffleNet-v2	ShuffleNet-v1
<i>Teacher</i>	74.64	79.34	79.34	79.42	79.42	75.61
<i>Student</i>	64.60	64.60	70.36	70.5	71.82	70.5
KD [11]	67.37	67.35	73.81	74.07	74.45	74.83
FitNet [22]	64.14 (\downarrow)	63.16 (\downarrow)	70.69 (\downarrow)	73.59 (\downarrow)	73.54 (\downarrow)	73.73 (\downarrow)
AT [29]	59.40 (\downarrow)	58.58 (\downarrow)	71.84 (\downarrow)	71.73 (\downarrow)	72.73 (\downarrow)	73.32 (\downarrow)
SP [26]	66.30 (\downarrow)	68.08 (\uparrow)	73.34 (\downarrow)	73.48 (\downarrow)	74.56 (\uparrow)	74.52 (\downarrow)
CC [21]	64.86 (\downarrow)	65.43 (\downarrow)	70.25 (\downarrow)	71.14 (\downarrow)	71.29 (\downarrow)	71.38 (\downarrow)
VID [1]	65.56 (\downarrow)	67.57 (\uparrow)	70.30 (\downarrow)	73.38 (\downarrow)	73.40 (\downarrow)	73.61 (\downarrow)
RKD [19]	64.52 (\downarrow)	64.43 (\downarrow)	71.50 (\downarrow)	72.28 (\downarrow)	73.21 (\downarrow)	72.21 (\downarrow)
PKT [20]	67.13 (\downarrow)	66.52 (\downarrow)	73.01 (\downarrow)	74.10 (\uparrow)	74.69 (\uparrow)	73.89 (\downarrow)
AB [10]	66.06 (\downarrow)	67.20 (\downarrow)	70.65 (\downarrow)	73.55 (\downarrow)	74.31 (\downarrow)	73.34 (\downarrow)
FT [13]	61.78 (\downarrow)	60.99 (\downarrow)	70.29 (\downarrow)	71.75 (\downarrow)	72.50 (\downarrow)	72.03 (\downarrow)
NST [12]	58.16 (\downarrow)	64.96 (\downarrow)	71.28 (\downarrow)	74.12 (\uparrow)	74.68 (\uparrow)	76.09 (\uparrow)
CRD [25]	69.73 (\uparrow)	69.11 (\uparrow)	74.3 (\uparrow)	75.11 (\uparrow)	75.65 (\uparrow)	76.05 (\uparrow)
CRD+KD [25]	69.94 (\uparrow)	69.54 (\uparrow)	74.58 (\uparrow)	75.12 (\uparrow)	76.05 (\uparrow)	76.27 (\uparrow)
OFD [9]	69.48 (\uparrow)	69.04 (\uparrow)	n/a	75.98 (\uparrow)	76.82 (\uparrow)	75.85 (\uparrow)
WSLD [32]	n/a	68.79 (\uparrow)	73.80 (\downarrow)	75.09 (\uparrow)	n/a	75.23 (\uparrow)
IPWD [18]	n/a	70.25 (\uparrow)	74.97 (\uparrow)	76.03 (\uparrow)	n/a	76.44 (\uparrow)
SRRL [27]	n/a	n/a	n/a	75.18 (\uparrow)	n/a	n/a
SemCKD [2]	n/a	n/a	n/a	n/a	77.62 (\uparrow)	n/a
ReviewKD [5]	70.37 (\uparrow)	69.89 (\uparrow)	n/a	77.45 (\uparrow)	77.78 (\uparrow)	77.14 (\uparrow)
SimKD [3]	n/a	n/a	n/a	77.18 (\uparrow)	n/a	n/a
DistPro [6]	n/a	n/a	n/a	77.18 (\uparrow)	77.54 (\uparrow)	77.24 (\uparrow)
NORM [14]	68.94 (\uparrow)	70.56 (\uparrow)	75.17 (\uparrow)	77.42 (\uparrow)	78.07 (\uparrow)	77.06 (\uparrow)
NORM+KD [14]	69.38 (\uparrow)	71.17 (\uparrow)	75.67 (\uparrow)	77.79 (\uparrow)	78.32 (\uparrow)	77.63 (\uparrow)
NORM+CRD [14]	69.17 (\uparrow)	71.08 (\uparrow)	75.51 (\uparrow)	77.50 (\uparrow)	77.96 (\uparrow)	77.09 (\uparrow)
WCoRD [4]	69.47 (\uparrow)	70.45 (\uparrow)	74.86 (\uparrow)	75.40 (\uparrow)	75.96 (\uparrow)	76.32 (\uparrow)
WCoRD+KD [4]	70.02 (\uparrow)	70.12 (\uparrow)	74.68 (\uparrow)	75.77 (\uparrow)	76.48 (\uparrow)	76.68 (\uparrow)
CRCD [33]	n/a	n/a	n/a	n/a	n/a	n/a
FKD [7]	n/a	n/a	74.61 (\uparrow)	75 (\uparrow)	n/a	n/a
ITRD (corr) [17]	69.97 (\uparrow)	71.41 (\uparrow)	75.71 (\uparrow)	76.8 (\uparrow)	77.27 (\uparrow)	77.35 (\uparrow)
ITRD (corr+mi) [17]	70.39 (\uparrow)	71.34 (\uparrow)	75.49 (\uparrow)	76.91 (\uparrow)	77.40 (\uparrow)	77.09 (\uparrow)
RRD (ours)	67.93 (\uparrow)	68.84 (\uparrow)	74.01 (\uparrow)	74.11 (\uparrow)	74.80 (\uparrow)	74.98 (\uparrow)
RRD+KD (ours)	69.98 (\uparrow)	69.13 (\uparrow)	74.26 (\uparrow)	75.18 (\uparrow)	76.29 (\uparrow)	76.31 (\uparrow)

conference on computer vision and pattern recognition, pages 4510–4520, 2018. 1

- [24] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2015. 1
- [25] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation, 2022. 1, 3, 4
- [26] Frederick Tung and Greg Mori. Similarity-preserving knowledge distillation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1365–1374, 2019. 1, 3, 4
- [27] Jing Yang, Brais Martinez, Adrian Bulat, and Georgios Tzimiropoulos. Knowledge distillation via softmax regres-

sion representation learning. In *International Conference on Learning Representations*, 2021. 1, 3, 4

- [28] 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*, pages 4133–4141, 2017. 1, 3
- [29] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *Proceedings of the 5th International Conference on Learning Representations*, 2016. 1, 3, 4
- [30] Sergey Zagoruyko and Nikos Komodakis. Wide residual

networks, 2017. 1

- [31] 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*, pages 6848–6856, 2018. 1
- [32] Helong Zhou, Liangchen Song, Jiajie Chen, Ye Zhou, Guoli Wang, Junsong Yuan, and Qian Zhang. Rethinking soft labels for knowledge distillation: A bias-variance tradeoff perspective, 2021. 1, 3, 4
- [33] Jinguo Zhu, Shixiang Tang, Dapeng Chen, Shijie Yu, Yakun Liu, Aijun Yang, Mingzhe Rong, and Xiaohua Wang. Complementary relation contrastive distillation, 2021. 1, 3, 4