

# Knowledge Distillation with the Reused Teacher Classifier

## – Supplementary Material –

Defang Chen<sup>1,2,3</sup> Jian-Ping Mei<sup>4</sup> Hailin Zhang<sup>1,2,3</sup>  
 Can Wang<sup>1,2,3</sup> Yan Feng<sup>1,2,3</sup> Chun Chen<sup>1,2,3</sup>

<sup>1</sup>Zhejiang University <sup>2</sup>Shanghai Institute for Advanced Study of Zhejiang University

<sup>3</sup>ZJU-Bangsun Joint Research Center <sup>4</sup>Zhejiang University of Technology

defchern@zju.edu.cn, jpmei@zjut.edu.cn, {zzzh1, wcan, fengyan, chenc}@zju.edu.cn

## A. Experimental Setting

### A.1. Datasets and Training Details

We adopt two datasets including CIFAR-100 [8] and ImageNet [13] to conduct experiments. All images are normalized by channel means and standard deviations. A horizontal flip is used for data augmentation. **CIFAR-100**<sup>1</sup> contains 50,000 training images and 10,000 test images from 100 classes. Each training image is padded by 4 pixels on each size and randomly cropped as a  $32 \times 32$  sample. **ImageNet**<sup>2</sup> contains about 1.3 million training images and 50,000 validation images from 1,000 classes. Each image is randomly cropped as a  $224 \times 224$  sample without padding. The top-1 test accuracy of the teacher model (ResNet-50) is 76.26%.

**Multi-Teacher Knowledge Distillation.** The training hyper-parameters of multi-teacher KD are exactly the same as those of single-teacher KD on CIFAR-100. We first pre-train multiple teacher models with different initialization and then distill their knowledge into a student model. The accuracies of compared AEKD and AEKD-F [4] are obtained by running a public library<sup>3</sup> with default model hyper-parameters on our teacher-student combinations [22]. The top-1 test accuracy of two groups of teacher models used in our main submission are: ① Three ResNet-32x4 models (79.32, 79.43, 79.45), ② Two ResNet-32x4 models (79.43, 79.45) and one ResNet-110x2 model (78.18).

**Data-Free Knowledge Distillation.** We adopt a public library<sup>4</sup> to reproduce the results of compared approaches: ZSKT [10], DAFL [3] and CMI [5], with the default model hyper-parameters. In our experiment, the top-1 test accuracy of the teacher model (WRN-40-2) is 76.31%. The performance of the student model trained with original dataset is included for comparison.

Input dimension	Operator	Output dimension
$H \times W \times C_s$	1x1 Conv	$H \times W \times C_t/r$
$H \times W \times C_t/r$	3x3 Conv	$H \times W \times C_t/r$
$H \times W \times C_t/r$	1x1 Conv	$H \times W \times C_t$

Table S.1. Projector structure. “1x1/3x3Conv” denotes a convolutional layer with 1x1/3x3 kernel size. Standard batch normalization and ReLU activation are used after each convolutional layer.  $r$  is the reduction ratio.

**Computing Infrastructure.** All of the experiments are conducted with PyTorch [11]. CIFAR-100 experiments are conducted on a sever containing eight NVIDIA GeForce RTX 2080Ti GPUs with 11GB RAM. The CUDA version is 11.2. ImageNet experiments are conducted on a sever containing four NVIDIA A40 GPUs with 48GB RAM. The CUDA version is 11.4.

### A.2. Network Architectures

We use a large number of teacher-student combinations for performance evaluation, which are composed of several popular neural network architectures: VGG [15], ResNet [6], WRN [20], MobileNetV2 [14], ShuffleNetV1 [23], ShuffleNetV2 [9]. The number behind “VGG-”, “ResNet-” denotes the depth of networks. “WRN-d-w” denotes the wide-ResNet with depth  $d$  and width factor  $w$ . As the previous works do [2, 16], we expand or shrink the number of convolution filters in intermediate layers of some network architectures with a certain ratio and put that ratio behind “x”, such as “ResNet-32x4”.

### A.3. Projector

The detailed structure of our used projector is described in Table S.1. We assume that the spatial dimensions of involved feature maps are the same, and denote them with the notations  $H$  and  $W$ . Otherwise, an average pooling op-

<sup>1</sup><https://www.cs.toronto.edu/~kriz/cifar.html>

<sup>2</sup><http://image-net.org/challenges/LSVRC/2012/index>

<sup>3</sup><https://github.com/Rorozhl/CA-MKD>

<sup>4</sup><https://github.com/zju-vipa/DataFree>

eration is used in advance for spatial dimension alignment to reduce the computational consumption, as the previous work do [2].

Given the feature maps of teacher and student models, the parameter number of the added projector is a function of the dimension reduction factor  $r$

$$\mathcal{F}(r) = \frac{C_t(C_s + C_t + 4)}{r} + \frac{9C_t^2}{r^2} + 2C_t. \quad (1)$$

**Proposition.** The extra parameter number  $\mathcal{F}(r)$  satisfies the inequality  $2\mathcal{F}(2r) < \mathcal{F}(r) < 4\mathcal{F}(2r)$  under some mild conditions.

**Proof.**

We first prove the left part of the inequality:

$$\begin{aligned} 2\mathcal{F}(2r) &< \mathcal{F}(r) \\ 2 \times \left( \frac{C_t(C_s + C_t + 4)}{2r} + \frac{9C_t^2}{4r^2} + 2C_t \right) &< \\ \frac{C_t(C_s + C_t + 4)}{r} + \frac{9C_t^2}{r^2} + 2C_t & \\ 2 \times \left( \frac{9C_t^2}{4r^2} + 2C_t \right) &< \frac{9C_t^2}{r^2} + 2C_t \\ 2C_t - \frac{9C_t^2}{2r^2} &< 0 \\ 2C_t \left( 1 - \frac{9C_t}{4r^2} \right) &< 0. \end{aligned} \quad (2)$$

Generally, the channel dimension  $C_t$  in the last feature maps of popular deep neural networks is greater than 128 on CIFAR-100 and is greater than 512 on ImageNet, which means that this equation holds when  $r < 16$  and  $r < 32$ , respectively. This is easy to be satisfied in practice. Since a typical setting for  $r$  is 1, 2 and 4 in order to avoid substantial accuracy reduction as shown in Table S.8 and S.9.

We then prove the right part of the inequality:

$$\begin{aligned} 4\mathcal{F}(2r) &> \mathcal{F}(r) \\ 4 \times \left( \frac{C_t(C_s + C_t + 4)}{2r} + \frac{9C_t^2}{4r^2} + 2C_t \right) &> \\ \frac{C_t(C_s + C_t + 4)}{r} + \frac{9C_t^2}{r^2} + 2C_t & \\ 4 \times \left( \frac{C_t(C_s + C_t + 4)}{2r} + 2C_t \right) &> \\ \frac{C_t(C_s + C_t + 4)}{r} + 2C_t & \\ 2 \times \left( \frac{C_t(C_s + C_t + 4)}{2r} + 3C_t \right) &> 0. \end{aligned} \quad (3)$$

Since the channel dimensions  $C_t$  and  $C_s$  are always greater than zero, this inequality holds automatically.  $\square$

Student	VGG-8	WRN-16-2	WRN-16-4
	70.46 $\pm$ 0.29	73.51 $\pm$ 0.32	77.26 $\pm$ 0.24
KD [7]	73.38 $\pm$ 0.05	75.40 $\pm$ 0.34	79.24 $\pm$ 0.23
FitNet [12]	73.63 $\pm$ 0.11	75.44 $\pm$ 0.22	79.06 $\pm$ 0.16
AT [21]	73.51 $\pm$ 0.08	75.76 $\pm$ 0.29	79.38 $\pm$ 0.20
SP [17]	73.53 $\pm$ 0.23	75.61 $\pm$ 0.34	79.53 $\pm$ 0.20
VID [1]	73.63 $\pm$ 0.07	75.44 $\pm$ 0.24	79.40 $\pm$ 0.08
CRD [16]	74.31 $\pm$ 0.17	75.86 $\pm$ 0.17	79.46 $\pm$ 0.19
SRRL [19]	74.25 $\pm$ 0.35	75.89 $\pm$ 0.12	79.67 $\pm$ 0.17
SemCKD [2]	74.43 $\pm$ 0.25	75.77 $\pm$ 0.11	80.05 $\pm$ 0.27
SimKD	<b>74.93 <math>\pm</math> 0.21</b>	<b>76.23 <math>\pm</math> 0.14</b>	<b>80.36 <math>\pm</math> 0.04</b>
Teacher	VGG-13	WRN-40-2	WRN-40-4
	74.64	76.31	79.51

Table S.2. Top-1 test accuracy (%) comparison on CIFAR-100.

Network	Student	SimKD	Teacher
ResNet-34 & ResNet-50	74.01	<b>74.64</b>	76.26
ResNet-50 & ResNet-101	76.26	<b>77.60</b>	77.80

Table S.3. Top-1 test accuracy (%) comparison on ImageNet.

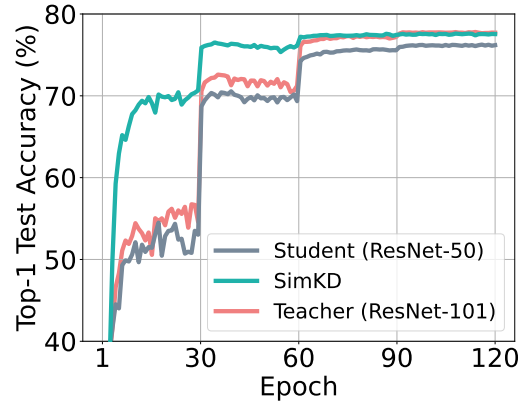


Figure S.1. The test accuracy (%) of ResNet-50 & ResNet-101 on ImageNet. Our SimKD achieves faster model convergence.

## B. More Experimental Results

### B.1. Comparison of Test Accuracy

Table S.2 and S.3 presents more results on CIFAR-100 and ImageNet datasets with extra *five* teacher-student combinations. Similar observations are obtained as those in the main submission. For ImageNet dataset, we replace the 3x3 convolution as the 3x3 depth-wise separable convolution in the projector (Table S.1) to control the extra parameters.

As shown in Figure S.1, our SimKD achieves *faster convergence* in the whole model training. For example, at 30th epoch, SimKD performance is on the par with the baseline student model performance at 60th epoch. Besides, at 60th epoch, SimKD already outperforms the baseline student model at 120th epoch.

Student	ResNet-8x4 73.09 $\pm$ 0.30		ResNet-8x4 73.09 $\pm$ 0.30	
	Student Classifier	Teacher Classifier	Student Classifier	Teacher Classifier
$\alpha = 0$ (KD)	<b>74.42 <math>\pm</math> 0.05</b>	–	<b>75.28 <math>\pm</math> 0.18</b>	–
$\alpha = 0.2$	74.42 $\pm$ 0.10	73.91 $\pm$ 0.12	74.83 $\pm$ 0.29	73.45 $\pm$ 0.31
$\alpha = 0.4$	73.99 $\pm$ 0.03	73.83 $\pm$ 0.13	74.72 $\pm$ 0.17	73.65 $\pm$ 0.27
$\alpha = 0.6$	73.93 $\pm$ 0.08	73.89 $\pm$ 0.26	74.59 $\pm$ 0.23	73.81 $\pm$ 0.25
$\alpha = 0.8$	73.76 $\pm$ 0.26	74.13 $\pm$ 0.24	74.29 $\pm$ 0.24	74.19 $\pm$ 0.22
$\alpha = 0.9$	73.52 $\pm$ 0.33	74.42 $\pm$ 0.26	73.98 $\pm$ 0.06	74.73 $\pm$ 0.13
$\alpha = 0.99$	30.77 $\pm$ 0.92	77.66 $\pm$ 0.22	23.46 $\pm$ 0.81	76.68 $\pm$ 0.13
$\alpha = 0.999$	5.59 $\pm$ 0.60	<b>77.98 <math>\pm</math> 0.19</b>	4.55 $\pm$ 0.68	<b>76.84 <math>\pm</math> 0.28</b>
$\alpha = 1$ (SimKD)	–	<b>78.08 <math>\pm</math> 0.15</b>	–	<b>76.75 <math>\pm</math> 0.23</b>
Teacher	ResNet-32x4 79.42		WRN-40-2 76.31	

Table S.4. Joint training the student feature encoder and classifier with different hyper-parameters.

Student	WRN-40-1 71.92 $\pm$ 0.17		MobileNetV2x2 69.06 $\pm$ 0.10	
	Student Classifier	Teacher Classifier	Student Classifier	Teacher Classifier
$\alpha = 0$ (KD)	<b>74.12 <math>\pm</math> 0.29</b>	–	<b>72.43 <math>\pm</math> 0.32</b>	–
$\alpha = 0.2$	74.33 $\pm$ 0.32	73.92 $\pm$ 0.31	72.50 $\pm$ 0.26	72.04 $\pm$ 0.23
$\alpha = 0.4$	74.49 $\pm$ 0.30	74.17 $\pm$ 0.19	72.48 $\pm$ 0.38	72.13 $\pm$ 0.25
$\alpha = 0.6$	74.53 $\pm$ 0.11	74.39 $\pm$ 0.20	72.83 $\pm$ 0.34	72.64 $\pm$ 0.29
$\alpha = 0.8$	74.93 $\pm$ 0.25	74.88 $\pm$ 0.22	73.13 $\pm$ 0.15	72.95 $\pm$ 0.06
$\alpha = 0.9$	75.18 $\pm$ 0.20	75.17 $\pm$ 0.17	73.44 $\pm$ 0.26	73.40 $\pm$ 0.27
$\alpha = 0.99$	74.17 $\pm$ 0.06	75.35 $\pm$ 0.15	74.91 $\pm$ 0.30	75.31 $\pm$ 0.16
$\alpha = 0.999$	18.39 $\pm$ 1.45	75.33 $\pm$ 0.08	13.49 $\pm$ 1.80	<b>75.43 <math>\pm</math> 0.22</b>
$\alpha = 1$ (SimKD)	–	<b>75.56 <math>\pm</math> 0.27</b>	–	<b>75.43 <math>\pm</math> 0.26</b>
Teacher	WRN-40-2 76.31		ResNet-32x4 79.42	

Table S.5. Joint training the student feature encoder and classifier with different hyper-parameters.

## B.2. Joint Training Results

Table S.4 and S.5 present the full joint training results with different hyper-parameters. In the case of  $\alpha = 0$  or  $\alpha = 1$ , only the student classifier or teacher classifier produces meaningful results and the another one degrades into random guess. We denote these random guess as “–”.

## B.3. Sequential Training Results

The results of sequential training in the main submission are obtained with the regular training procedure. That is to say, we adopt SGD with 0.9 Nesterov momentum and  $5 \times 10^{-4}$  weight decay. The total training epoch is set to 240 and the learning rate is divided by 10 at 150th, 180th

Learning Rate	Test Accuracy
0.01	52.03 $\pm$ 0.15
<b>0.05</b>	<b>51.97 <math>\pm</math> 0.19</b>
0.1	52.01 $\pm$ 0.17
0.5	51.93 $\pm$ 0.20

Table S.6. Training a new classifier from scratch with different initial learning rates (Student: ResNet-8x4, Teacher: ResNet-32x4).

and 210th epochs. The initial learning rate is set to 0.01 for MobileNet/ShuffleNet-series architecture and 0.05 for other architectures. The mini-batch size is set to 64.

	Input ( $\ell_2$ )	Output ( $\ell_2$ )	Input ( $\ell_2$ ) + Output ( $\ell_2$ )
Accuracy	<b>78.08 <math>\pm</math> 0.15</b>	77.09 $\pm$ 0.09	77.88 $\pm$ 0.30
Loss function	$\ \mathbf{f}^t - \mathbf{f}^s\ _2^2$	$\ \mathbf{W}^t \mathbf{f}^t - \mathbf{W}^t \mathbf{f}^s\ _2^2$	$\ \mathbf{f}^t - \mathbf{f}^s\ _2^2 + \ \mathbf{W}^t \mathbf{f}^t - \mathbf{W}^t \mathbf{f}^s\ _2^2$
Gradient on $\mathbf{f}^s$	$-2(\mathbf{f}^t - \mathbf{f}^s)$	$-2\mathbf{W}^{tT}\mathbf{W}^t(\mathbf{f}^t - \mathbf{f}^s)$	$-2\{(\mathbf{I} + \mathbf{W}^{tT}\mathbf{W}^t)(\mathbf{f}^t - \mathbf{f}^s)\}$

Table S.7. Comparison of different loss functions (Student: ResNet-8x4, Teacher: ResNet-32x4). We omit the projector  $\mathcal{P}(\cdot)$  for simplicity.

Student	WRN-40-1 71.92 $\pm$ 0.17	ResNet-8x4 73.09 $\pm$ 0.30	ResNet-110 74.37 $\pm$ 0.17	ResNet-116 74.46 $\pm$ 0.09	VGG-8 70.46 $\pm$ 0.29	ResNet-8x4 73.09 $\pm$ 0.30	ShuffleNetV2 72.60 $\pm$ 0.12
$r = 8$	67.20 $\pm$ 0.35 (0.55%, 1.05%)	<b>76.73 <math>\pm</math> 0.20</b> ( <b>0.35%</b> , 2.11%)	71.71 $\pm$ 1.00 (0.18%, 0.35%)	71.96 $\pm$ 1.09 (0.18%, 0.33%)	<b>74.74 <math>\pm</math> 0.15</b> ( <b>0.90%</b> , 0.86%)	66.26 $\pm$ 0.98 (-0.17%, 0.73%)	77.49 $\pm$ 0.31 (-0.35%, 3.76%)
$r = 4$	74.29 $\pm$ 0.03 (0.99%, 2.81%)	<b>77.88 <math>\pm</math> 0.41</b> ( <b>0.94%</b> , 5.67%)	<b>77.14 <math>\pm</math> 0.22</b> ( <b>0.32%</b> , 0.92%)	<b>77.18 <math>\pm</math> 0.21</b> ( <b>0.32%</b> , 0.87%)	<b>75.62 <math>\pm</math> 0.28</b> ( <b>1.62%</b> , 2.19%)	75.57 $\pm$ 0.03 (0.41%, 1.78%)	<b>78.21 <math>\pm</math> 0.20</b> ( <b>0.58%</b> , 8.85%)
$r = 2$	<b>75.56 <math>\pm</math> 0.27</b> ( <b>2.5%</b> , 8.77%)	<b>78.08 <math>\pm</math> 0.15</b> ( <b>2.88%</b> , 17.34%)	<b>77.82 <math>\pm</math> 0.15</b> ( <b>0.82%</b> , 2.88%)	<b>77.90 <math>\pm</math> 0.11</b> ( <b>0.82%</b> , 2.73%)	<b>75.76 <math>\pm</math> 0.12</b> ( <b>2.98%</b> , 6.23%)	<b>76.75 <math>\pm</math> 0.23</b> ( <b>2.18%</b> , 5.02%)	<b>78.39 <math>\pm</math> 0.27</b> ( <b>3.16%</b> , 23.01%)
$r = 1$	75.95 $\pm$ 0.30 (7.95%, 30.35%)	78.80 $\pm$ 0.13 (9.71%, 58.51%)	<b>78.00 <math>\pm</math> 0.26</b> ( <b>2.6%</b> , 9.96%)	<b>78.15 <math>\pm</math> 0.30</b> ( <b>2.6%</b> , 9.43%)	75.98 $\pm$ 0.21 (11.05%, 19.87%)	76.96 $\pm$ 0.07 (8.16%, 15.96%)	78.66 $\pm$ 0.08 (11.33%, 67.77%)
Teacher	WRN-40-2 76.31	ResNet-32x4 79.42	ResNet-110x2 78.18	ResNet-110x2 78.18	ResNet-32x4 79.42	WRN-40-2 76.31	ResNet-32x4 79.42

Table S.8. Top-1 test accuracy (%) and pruning ratio (the first element in parenthesis) of SimKD with various dimension reduction factor  $r$  on CIFAR-100. We also provide the ratio of the projector parameters to the student parameters (the second element in parenthesis).

Table S.6 gives additional results with different initial learning rates. It is shown that the student accuracy always stays at about 50% when the learning rate ranges from 0.01 to 0.5, which indicates the difficulty of training a satisfactory student classifier from scratch. In contrast, our SimKD achieves **78.08  $\pm$  0.15** test accuracy without any classifier retraining but just reusing the pre-trained teacher classifier.

#### B.4. Comparison of Loss Function

The default feature alignment loss in our main submission is implemented in the preceding layer of the teacher classifier with a  $\ell_2$  loss. Its result is reported in the second column of Table S.7. Another implementation is to calculate the loss in the succeeding layer of the teacher classifier with a loss function  $\|\mathbf{W}^t \mathbf{f}^t - \mathbf{W}^t \mathcal{P}(\mathbf{f}^s)\|_2^2$ , and we report its results in the third column of Table S.7. From Table S.7, we find that our default feature alignment loss performs best. Moreover, the gradient comparison of different loss functions indicates that the effect of “Output ( $\ell_2$ )” is to calibrate the gradient of “Input ( $\ell_2$ )” with a symmetric matrix  $\mathbf{W}^{tT}\mathbf{W}^t$ .

#### B.5. Comparison of Pruning Ratio

Table S.8 and S.9 present the top-1 test accuracy and the cost of pruning ratio (the first element in parenthesis) of SimKD versus different dimension reduction factors. We also provide the ratio of the projector parameters to the student parameters (the second element in parenthesis) for

comparison. We make those results bold when SimKD achieves state-of-the-art performance and the added projector only requires less than or about 3% pruning ratio cost.

In some cases such as “MobileNetV2x2 & ResNet-32x4” and “ShuffleNetV1 & ResNet-32x4” with  $r = 8$ , we can see that the pruning ratios of SimKD are even higher than the vanilla KD training, and all competitors accordingly. Moreover, SimKD achieves the second best performance on “ShuffleNetV2 & ResNet-32x4” with  $r = 8$  (SimKD: 77.49%, the best performance is achieved by SemCKD: 77.62%), “ShuffleNetV2x1.5 & ResNet-32x4” with  $r = 8$  (SimKD: 78.96%, the best performance is achieved by SemCKD: 79.13%), and “ShuffleNetV2 & ResNet-110x2” with  $r = 4$ , (SimKD: 77.35%, the best performance is achieved by SemCKD: 77.67%). Although the projector needs retaining during the whole training and test stages, a series of trade-off experiments between test accuracy and pruning ratio show that the extra parameters it brought are negligible in most cases.

We further extend our technique to the situation where more deep teacher layers are reused for student inference and analyze the accompanying trade-off between accuracy enhancement and complexity increase. As shown in Table S.10, “SimKD+” and “SimKD++” achieve higher performance than “SimKD” but they also bring about a sharp drop of the pruning ratio, which indicates that simply reusing the final teacher classifier strikes a good balance between performance and parameter complexity.

Student	ShuffleNetV1 71.36 $\pm$ 0.25	WRN-16-2 73.51 $\pm$ 0.32	ShuffleNetV2 72.60 $\pm$ 0.12	MobileNetV2 65.43 $\pm$ 0.29	MobileNetV2x2 69.06 $\pm$ 0.10	WRN-40-2 76.35 $\pm$ 0.18	ShuffleNetV2x1.5 74.15 $\pm$ 0.22
$r = 8$	<b>76.68 <math>\pm</math> 0.20</b> (-0.29%, 5.16%)	75.41 $\pm$ 0.17 (0.47%, 3.13%)	73.50 $\pm$ 0.77 (-0.99%, 1.55%)	61.78 $\pm$ 1.21 (-4.00%, 3.08%)	<b>74.97 <math>\pm</math> 0.14</b> (-0.58%, 2.51%)	78.55 $\pm$ 0.33 (0.13%, 0.98%)	78.96 $\pm$ 0.15 (-0.35%, 1.98%)
$r = 4$	<b>77.22 <math>\pm</math> 0.22</b> (0.60%, 12.12%)	<b>76.69 <math>\pm</math> 0.23</b> (1.01%, 8.81%)	77.35 $\pm$ 0.18 (-0.63%, 3.39%)	69.43 $\pm$ 0.21 (-2.67%, 6.77%)	<b>75.56 <math>\pm</math> 0.34</b> (0.45%, 5.78%)	<b>79.23 <math>\pm</math> 0.06</b> (0.27%, 2.75%)	<b>79.48 <math>\pm</math> 0.12</b> (0.58%, 4.65%)
$r = 2$	<b>77.18 <math>\pm</math> 0.26</b> (3.14%, 32.03%)	<b>77.17 <math>\pm</math> 0.32</b> (2.84%, 28.13%)	<b>78.25 <math>\pm</math> 0.24</b> (0.31%, 8.19%)	<b>70.71 <math>\pm</math> 0.41</b> (0.52%, 15.62%)	<b>75.43 <math>\pm</math> 0.26</b> (3.26%, 14.66%)	<b>79.29 <math>\pm</math> 0.11</b> (0.76%, 8.78%)	<b>79.54 <math>\pm</math> 0.26</b> (3.16%, 12.09%)
$r = 1$	77.58 $\pm$ 0.36 (11.20%, 95.15%)	77.65 $\pm$ 0.24 (9.45%, 98.01%)	78.58 $\pm$ 0.22 (2.99%, 21.83%)	70.90 $\pm$ 0.17 (9.43%, 40.34%)	75.71 $\pm$ 0.20 (11.87%, 41.86%)	79.26 $\pm$ 0.17 (2.55%, 30.59%)	79.72 $\pm$ 0.24 (11.33%, 35.61%)
Teacher	ResNet-32x4 79.42	ResNet-32x4 79.42	ResNet-110x2 78.18	WRN-40-2 76.31	ResNet-32x4 79.42	ResNet-110x4 80.20	ResNet-32x4 79.42

Table S.9. Top-1 test accuracy (%) and pruning ratio (the first element in parenthesis) of SimKD with various dimension reduction factor  $r$  on CIFAR-100. We also provide the ratio of the projector parameters to the student parameters (the second element in parenthesis).

	Test Accuracy	Pruning Ratio
Student	73.09 $\pm$ 0.30	83.40%
SimKD	<b>78.08 <math>\pm</math> 0.15</b>	<b>80.52%</b>
SimKD+	<b>78.47 <math>\pm</math> 0.08</b>	19.21%
SimKD++	<b>78.88 <math>\pm</math> 0.05</b>	15.97%
Teacher	79.42	0%

Table S.10. Comparison of reusing different teacher layers.

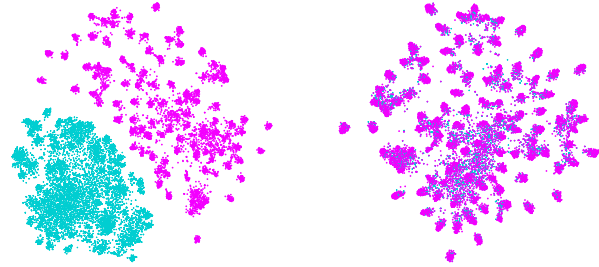
## B.6. Visualization

We adopt ResNet-8x4 as the student model and ResNet-32x4 as the teacher model for visualization experiments.

Ten randomly selected classes in the main submission includes “road”, “bee”, “lawn\_mower”, “bottle”, “shrew”, “bridge”, “man”, “mouse”, “sweet\_pepper” and “cattle”. We further visualize all 100 classes on CIFAR-100 with t-SNE in Figure S.2. The visualization results show that with the help of a simple  $\ell_2$  loss, the extracted features from teacher and student models become almost indistinguishable in SimKD, which ensures the student features to be correctly classified with the reused teacher classifier later.

## References

- [1] Sungsoo Ahn, Shell Xu Hu, Andreas C. 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. 2
- [2] Defang Chen, Jian-Ping Mei, Yuan Zhang, Can Wang, Zhe Wang, Yan Feng, and Chun Chen. Cross-layer distillation with semantic calibration. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 7028–7036, 2021. 1, 2
- [3] Hanting Chen, Yunhe Wang, Chang Xu, Zhaohui Yang, Chuanjian Liu, Boxin Shi, Chunjing Xu, Chao Xu, and Qi Tian. Data-free learning of student networks. In *International Conference on Computer Vision*, pages 3513–3521, 2019. 1



(a) Vanilla KD [7].

(b) Our SimKD.

Figure S.2. Visualizations of all test images from CIFAR-100 with t-SNE [18]. Features extracted by the teacher and student models are depicted with magenta and cyan colors, respectively, and they are almost indistinguishable in our SimKD. Best viewed in color.

- [4] Shangchen Du, Shan You, Xiaojie Li, Jianlong Wu, Fei Wang, Chen Qian, and Changshui Zhang. Agree to disagree: Adaptive ensemble knowledge distillation in gradient space. In *Advances in Neural Information Processing Systems*, 2020. 1
- [5] Gongfan Fang, Jie Song, Xinchao Wang, Chengchao Shen, Xingen Wang, and Mingli Song. Contrastive model inversion for data-free knowledge distillation. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pages 2374–2380, 2021. 1
- [6] 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*, pages 770–778, 2016. 1
- [7] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network, 2015. 2, 5
- [8] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. *Technical Report*, 2009. 1
- [9] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet V2: practical guidelines for efficient CNN architecture. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1798–1807, 2018. 1

- ecture design. In *Proceedings of the European Conference on Computer Vision*, pages 122–138, 2018. 1
- [10] Paul Micaelli and Amos J. Storkey. Zero-shot knowledge transfer via adversarial belief matching. In *Advances in Neural Information Processing Systems*, pages 9547–9557, 2019. 1
- [11] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, pages 8024–8035, 2019. 1
- [12] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fit-nets: Hints for thin deep nets. In *International Conference on Learning Representations*, 2015. 2
- [13] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Fei-Fei Li. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. 1
- [14] Mark Sandler, Andrew G. 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*, pages 4510–4520, 2018. 1
- [15] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015. 1
- [16] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. In *International Conference on Learning Representations*, 2020. 1, 2
- [17] Frederick Tung and Greg Mori. Similarity-preserving knowledge distillation. In *International Conference on Computer Vision*, pages 1365–1374, 2019. 2
- [18] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(86):2579–2605, 2008. 5
- [19] Jing Yang, Brais Martínez, Adrian Bulat, and Georgios Tzimiropoulos. Knowledge distillation via softmax regression representation learning. In *International Conference on Learning Representations*, 2021. 2
- [20] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *Proceedings of the British Machine Vision Conference*, 2016. 1
- [21] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: improving the performance of convolutional neural networks via attention transfer. In *International Conference on Learning Representations*, 2017. 2
- [22] Hailin Zhang, Defang Chen, and Can Wang. Confidence-aware multi-teacher knowledge distillation. *arXiv preprint arXiv:2201.00007*, 2021. 1
- [23] 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