

# Semi-Online Knowledge Distillation

Zhiqiang Liu<sup>1</sup>  
sezhiqiangliu@mail.scut.edu.cn  
Yanxia Liu<sup>\*,1</sup>  
cslyx@scut.edu.cn  
Chengkai Huang<sup>2</sup>  
fjshwhck@hrbeu.edu.cn

<sup>1</sup> School of Software  
South China University of Technology  
Guangzhou, China  
<sup>2</sup> College of Computer Science and  
Technology  
Harbin Institute of Technology  
Shenzhen, China

## Abstract

Knowledge distillation is an effective and stable method for model compression via knowledge transfer. Conventional knowledge distillation (KD) is to transfer knowledge from a large and well pre-trained teacher network to a small student network, which is a one-way process. Recently, deep mutual learning (DML) has been proposed to help student networks learn collaboratively and simultaneously. However, to the best of our knowledge, KD and DML have never been jointly explored in a unified framework to solve the knowledge distillation problem. In this paper, we investigate that the teacher model supports more trustworthy supervision signals in KD, while the student captures more similar behaviors from the teacher in DML. Based on these observations, we first propose to combine KD with DML in a unified framework. Furthermore, we propose a Semi-Online Knowledge Distillation (SOKD) method that effectively improves the performance of the student and the teacher. In this method, we introduce the peer-teaching training fashion in DML in order to alleviate the student’s imitation difficulty, and also leverage the supervision signals provided by the well-trained teacher in KD. Besides, we also show our framework can be easily extended to feature-based distillation methods. Extensive experiments on CIFAR-100 and ImageNet datasets demonstrate the proposed method achieves state-of-the-art performance.

## 1 Introduction

In recent years, deep learning has demonstrated great success in many computer vision tasks such as image classification [10, 15, 24], object detection [18, 22, 33] and image generation [0, 19, 44], which mainly benefits from deeper and larger convolutional neural network architectures. However, these architectures require huge computation and massive storage, which makes them hard to deploy on resource-limited devices such as mobile phones and drones. Therefore, many recent studies have focused on network compression methods, including network pruning [6, 11, 45], network quantization [29, 30, 31], and knowledge distillation [12, 23, 35].

Knowledge distillation is a simple yet effective technology, which is orthogonal to other model compression methods. Existing distillation methods can be roughly divided into offline knowledge distillation and online knowledge distillation [9]. Without loss of generality,

in this paper we investigate the typical offline and online knowledge distillation methods, namely KD [14] and deep mutual learning (DML) [15]. The key idea of KD is to improve the performance of a compact network (student) via the knowledge transfer from a large and well-trained network (teacher). Specifically, Hinton *et al.* [14] regard the softened output probabilities of the teacher as knowledge and propose that the student absorbs this knowledge by approximating the teacher’s outputs. Whereas, the student is generally hard to fully absorb the knowledge provided by the teacher in practice (i.e., fully match the outputs of the teacher). The main reason is that there exists a large performance gap between the converged large network and untrained small network [16]. On the contrary, the networks in DML are easier to approach and learn from each other because they are regarded as students and trained simultaneously. Such peer-teaching training fashion makes the networks easily imitate the behaviors (i.e., representations and outputs) of each other. However, since all networks in DML are trained from scratch and their optimization directions seem unstable in the early training phase, there exist many unreliable supervision signals, which might lead to conflicts of networks’ optimization [9]. Although both KD and DML exist some shortcomings, we argue that they are complementary to each other for knowledge distillation problem.

To verify the above viewpoint, we evaluate the imitation ability of the student<sup>1</sup> based on CKA similarity [17] and two well-defined metrics which are imitation error rate (IER) and misleading rate (MR). We empirically found that 1) The peer-teaching training fashion in DML is a more effective way than one-way flow in KD for the student to learn similar behaviors of the teacher. 2) The well-trained teacher in KD can provide more trustworthy and stable supervision signals for the student than the teacher in DML. (See discussion in detail in section 3.2).

Based on the above observations, in this paper we first propose a unified framework combining KD with DML, which effectively leverages the teacher’s supervision signals in KD and alleviates the imitation difficulty via the peer-teaching training fashion in DML. To this end, we further propose a Semi-Online Knowledge Distillation (SOKD) method. In this method, we build a simple yet effective knowledge bridge module (KBM) to exploit the knowledge from the offline teacher and transfer it to the student online. Specifically, the KBM is the same as the high-level layers of the teacher while taking the output of the low-level layers as input. In the training phase, inspired by DML, we update the KBM and the student simultaneously while fixing the teacher. In the inference phase, we replace high-level layers with KBM to construct a new teacher. As a result, we obtain a compact but effective student as well as, surprisingly, a more powerful teacher.

Our contributions are summarized as follows:

- We discover that KD and DML are complementary to each other on knowledge distillation task. Based on this observation, we propose a novel distillation method combining KD with DML. To the best of our knowledge, we are the first to jointly integrate KD with DML in a unified framework.
- We propose a Semi-Online Knowledge Distillation (SOKD) method that significantly improves the performance of the student and the teacher. To leverage the well-trained teacher’s supervision signals and alleviate the imitation difficulty of the student, we propose a simple yet effective knowledge bridge module (KBM).

<sup>1</sup>For convenience, we only consider two networks case and refer to the smaller network in DML as student, otherwise as teacher.

- We further extend our proposed framework to feature-based distillation methods. Extensive experiments on CIFAR-100 and ImageNet datasets demonstrate that the proposed SOKD not only outperforms state-of-the-art model distillation methods but also can be easily extended to feature-based distillation methods for improved learning accuracy.

## 2 Related Work

**Offline Knowledge Distillation.** The training process of traditional offline knowledge distillation can be split into two stages: 1) pre-training a teacher; and 2) distilling the knowledge of teacher into student. Most existing work mainly concentrate on how to define and transfer knowledge. For example, Hinton *et al.* [12] first propose to mimic softened logits of teacher. Romero *et al.* [23] propose to match the intermediate representations of the teacher and the student. After that, many studies investigate how to clearly define the knowledge based on the intermediate feature maps [27, 33]. However, most feature-based distillation methods are sensitive to the structures of the teacher and the student [26, 32]. Instead, in this paper we focus on improving the student’s capacity to imitate the logits of the teacher. Such logit-based knowledge distillation can be well generalized to different teacher-student combinations.

**Online Knowledge Distillation.** Zhang *et al.* [42] propose an online distillation scheme called deep mutual learning (DML), in which all networks are treated as students and updated simultaneously. To reduce the computational cost, Lan *et al.* [43] propose a multi-branch network. Recently, Chen *et al.* [9] introduce the self-attention mechanism to improve the diversity of students. However, the different outputs of students will conflict with each other, which may harm the convergence of training [9]. To improve the ability of converging with higher generalization, Guo *et al.* [9] generate a high-quality soft target via ensembling the output of students. Besides, some studies [13, 36, 40] propose self-distillation methods, the special cases of online knowledge distillation, in which the teachers and students are the same networks. In this paper, we introduce the peer-teaching training fashion in DML into KD to simplify the difficulty for student to imitate teacher’s behaviors.

## 3 Proposed Method

### 3.1 Preliminary

**Knowledge Distillation.** The key idea of KD is that the student uses the softened output of the teacher as supervised signal for network training. Given the  $i^{th}$  sample, we can get the logits  $z_i = (z_i^1, z_i^2, \dots, z_i^C)$  after softmax function. Then the softened output probability of the teacher network is computed as:

$$p_i^t = \frac{\exp(z_i^j / \tau)}{\sum_{j=1}^C \exp(z_i^j / \tau)}, \quad (1)$$

where  $\tau$  is a temperature factor. Similarly, the student produces a softened target  $p_i^s$ . The student mimics outputs of the teacher by minimizing Kullback-Leibler (KL) Divergence  $KL(p^s, p^t)$ . The total loss for the student network consists of the typical cross entropy loss

$\mathcal{L}_{ce}^s$  and the KL Divergence:

$$KL(p^s, p^t) = \sum_{x \sim \mathcal{D}} p^s \log \frac{p^s}{p^t}, \mathcal{L}^s = \lambda_1 \mathcal{L}_{ce}^s + \lambda_2 KL(p^s, p^t), \quad (2)$$

$\lambda_1$  and  $\lambda_2$  are hyperparameters.

**Deep Mutual Learning.** Zhang *et al.* [42] propose an online knowledge distillation method named deep mutual learning, in which the teacher and the student networks are updated simultaneously. Specifically, the teacher and the student optimize their networks with their respective loss functions  $\mathcal{L}^t$  and  $\mathcal{L}^s$ . The formulations are as follows:

$$\mathcal{L}^t = \lambda_1 \mathcal{L}_{ce}^t + \lambda_2 KL(p^t, p^s), \mathcal{L}^s = \lambda_1 \mathcal{L}_{ce}^s + \lambda_2 KL(p^s, p^t), \quad (3)$$

$\lambda_1$  and  $\lambda_2$  are hyperparameters.

### 3.2 Analysis on Imitation Ability of Student

In distillation methods, the student is expected to imitate the behaviors of the teacher, which means that the performance of the student largely depends on its imitation. To investigate the mimic ability of the student in KD and DML, we conduct experiments with the teacher-student pair of WRN-28-4 and WRN-16-2 on CIFAR-100 dataset. Specifically, we use CKA [43] to measure the representational similarity between the teacher and the student. The larger the CKA, the more similar the representations between teacher and student. Besides, we define the imitation error rate (IER) and the misleading rate (MR) as the metrics to measure the ability to imitate the final output. IER refers to the rate of different outputs between the teacher and the student. MR represents the rate of wrong outputs learned from the teacher. Let  $\mathcal{D}_{st}$  be the samples with the same prediction results between the teacher and the student in the training data  $\mathcal{D}$ ,  $\mathcal{D}_{sg}$  be the samples with the same prediction results between the ground-truth and the student in  $\mathcal{D}_{st}$ . IER and MR are computed as  $(1 - |\mathcal{D}_{st}|/|\mathcal{D}|) \times 100$  and  $(1 - |\mathcal{D}_{sg}|/|\mathcal{D}_{st}|) \times 100$ , respectively.  $|\cdot|$  denotes number of samples.

From table 5, DML achieves much higher CKA and lower IER than KD, which means that the student in DML has stronger abilities to capture the representations of the teacher and mimic the output probabilities of the teacher effectively. However, DML achieves highest MR, meaning that the student in DML has learned much more wrong information from the teacher. The main reason is that in the early training phase the optimization direction of the teacher is unstable. The supervision signals provided by the teacher might be inaccurate and changing. Fortunately, we are able to easily obtain reliable supervision signals from the well-trained teacher in KD. Therefore, it is possible for the student to learn similar and accurate behaviors from the teacher and effectively improve its performance.

### 3.3 Semi-Online Knowledge Distillation

Based on above analysis, the training mechanism that the teacher and the student are updated simultaneously in DML is an effective way for the student to learn similar behaviors of the teacher. Besides, the well-trained teacher in KD provides reliable and stable signals for the student. In this paper, we propose to integrate KD with DML in a unified framework. To this end, we propose a Semi-Online Knowledge Distillation (SOKD) framework that effectively improves the performance of the student and the teacher. The overview of our framework

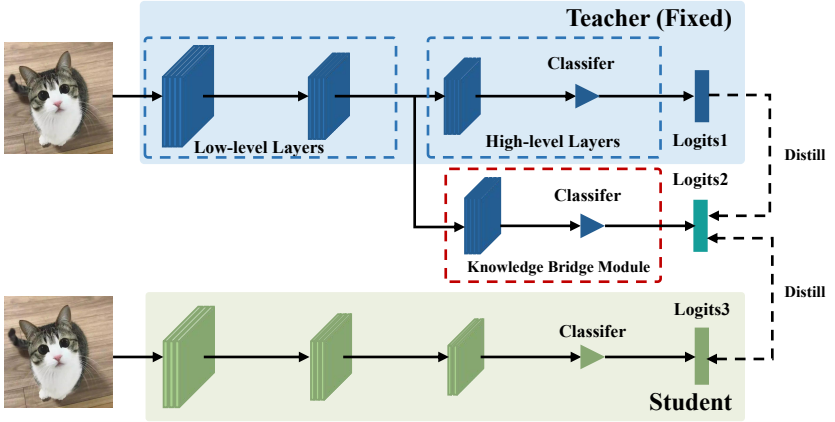


Figure 1: An overview of the proposed SOKD. In the training phase, the teacher is fixed. The knowledge bridge module (KBM) is trained under the supervision of both the teacher and the student. The student is to mimic the output of the KBM. In the inference phase, we reconstruct the teacher network with low-level layers and KBM. In this way, we obtain a well-trained student network as well as a retrained teacher network.

is shown in figure 1, which consists of three parts, well-trained teacher (low-level and high-level layers), knowledge bridge module (KBM), and student. The KBM is proposed to build a knowledge bridge between the teacher and the student. We take the outputs of low-level layers as the inputs of KBM. More details about KBM will be discussed in section 3.4.

In the training phase, the KBM tries to fully exploit the knowledge of the teacher and transfer it to the student. Meanwhile, the KBM also gains the feedback from the student. Specifically, we do not update the parameters of the teacher since it would destroy the inherent knowledge of the teacher. Instead, we update the parameters of the KBM and the student simultaneously, which is a more tolerant way for the student to approach what the teacher expresses. The training objectives of the KBM and the student can be written as:

$$\begin{aligned}\mathcal{L}^{kbn} &= \alpha_1 \mathcal{L}_{ce}^{kbn} + \alpha_2 KL(p^{kbn}, p^t) + \alpha_3 KL(p^{kbn}, p^s), \\ \mathcal{L}^s &= \lambda_1 \mathcal{L}_{ce}^s + \lambda_2 KL(p^s, p^{kbn}),\end{aligned}\quad (4)$$

where  $\alpha_1$ ,  $\alpha_2$ ,  $\alpha_3$ ,  $\lambda_1$ ,  $\lambda_2$  are balanced hyperparameters. KL divergence is computed by Eqn. 2. In this sense, we are able to obtain a compact but effective student network. Besides, we also get a more powerful teacher by reconstructing the teacher network with low-level layers and KBM.

### 3.4 Knowledge Bridge Module

In this subsection, we discuss how to build the KBM, the structure of which is motivated by characteristics of Convolutional Neural Network (CNN). CNN can be roughly composed of low-level and high-level layers. Given a face image, the features detected by low-level layers are usually some fundamental patterns, *e.g.*, edge, corner, and color, which are not sensitive to classes [2]. By contrast, the features extracted by high-level layers may be a full face, which is class-specific and related to the output closely. In other words, the discriminative power of a CNN is mainly reflected in the high-level layers [9].

Teacher Student	WRN-28-4 WRN-16-2	WRN-40-2 WRN-40-1	resnet110 resnet20	vgg13 vgg8	MBNV2 (1.4) MBNV2 (0.5)	ShuffleNetV1 MBNV2 (0.5)	resnet32×4 ShuffleNetV1	WRN-28-4 resnet56
Teacher	78.91	76.22	74.00	75.29	69.57	71.85	78.92	78.91
Student	73.32	71.78	69.20	70.99	65.17	65.17	71.85	73.05
KD [□]	74.72	73.54	70.73	73.58	67.94	68.19	73.90	74.21
AT [▣]	74.01	72.94	70.78	72.78	66.26	65.55	73.15	73.64
SP [▣]	73.81	73.18	70.69	73.58	61.52	67.67	74.87	74.28
CC [▣]	73.73	72.22	70.30	71.49	65.39	65.16	72.41	73.16
VID [□]	73.93	72.52	70.26	72.13	65.15	65.49	73.10	73.24
CRD [▣]	75.05	73.57	70.88	73.24	66.47	67.32	73.76	74.34
SSKD [▣]	75.57	74.02	70.70	73.96	68.37	67.99	73.99	75.31
SemCKD [□]	73.63	73.32	70.41	73.35	67.70	67.19	74.29	73.13
SOKD (Teacher)	80.54	78.41	76.32	77.16	71.83	73.41	80.46	80.47
SOKD (Student)	<b>76.82</b>	<b>75.35</b>	<b>72.08</b>	<b>74.27</b>	<b>69.28</b>	<b>69.50</b>	<b>75.87</b>	<b>75.94</b>

Table 1: Comparison with state-of-the-art offline methods on CIFAR-100 dataset.

Networks		Vanilla		DML [▣]		KDCL [□]		DCM [▣]		SOKD (Ours)	
Teacher	Student	Teacher	Student	Teacher	Student	Teacher	Student	Teacher	Student	Teacher	Student
WRN-16-2	WRN-16-2	73.32	73.32	74.46	74.56	74.64	74.57	74.88	74.86	75.54	<b>75.63</b>
WRN-28-2	WRN-28-2	75.75	75.75	76.83	77.00	77.36	77.32	77.71	77.67	78.04	<b>78.61</b>
WRN-28-4	WRN-16-2	78.91	73.32	79.69	75.27	79.89	75.14	78.92	75.44	<b>80.54</b>	<b>76.82</b>
WRN-28-4	MobileNet	78.91	73.23	80.49	76.94	80.64	76.03	<b>80.67</b>	78.11	80.62	<b>78.84</b>

Table 2: Comparison with state-of-the-art online methods on CIFAR-100 dataset. We refer to the smaller network in online methods as student, otherwise as teacher.

Based on the above theoretical analyses, we design the KBM to be the same as the teacher’s high-level layers while using its low-level layer’s output as input. Formally, given a pre-trained teacher model  $M^t$ , our main goal is to obtain a compact yet efficient student model  $M^s$ . For simplicity, low-level and high-level layers are denoted by  $\{M_i^t\}_{i=1}^l$  and  $\{M_i^t\}_{i=l+1}^L$ , respectively. To build a knowledge bridge between  $M^t$  and  $M^s$ , we propose the KBM  $M^{kbm}$ . The KBM is the same as  $\{M_i^t\}_{i=l+1}^L$  and takes the output of  $\{M_i^t\}_{i=1}^l$  as input.

The main characteristics of the KBM are: 1) Since KBM is the same as high-level layers of the teacher, it has the ability to approach what the teacher expresses. 2) Since different networks have similar representations in low-level layers, sharing the low-level layers would not greatly affect the performance of the KBM. Conversely, it is similar to layer-wise training [3, 23], which boosts training stability. Besides, it only brings a small amount of extra calculations.

### 3.5 Extension to Feature-Based Methods

In this subsection, we investigate the possibility of extending our method to feature-based methods. The main differences between feature-based and logit-based methods lie in the definition of the knowledge and learning objectives of the student. Feature-based methods learn from the teacher by minimizing the distance (*e.g.*,  $L_1$ ,  $L_2$  distance) of intermediate activations. Different methods have different requirements in the positions of intermediate activations. For example, SP only uses the activation of last convolutional layer, while AT adopts activations at the end of each feature extractor block (*e.g.*, residual block). However, it would not be an obstacle of the extension of our method. Our proposed KBM structure can be easily adjusted according to the intermediate activations without loss of performance. The KBM with different blocks will be discussed in detail in section 4.4. Formally, the overall

	Teacher	Student	KD [15]	AT [16]	SP [17]	CC [18]	CRD [19]	SSKD [20]	SemCKD [8] <sup>†</sup>	SOKD (Teacher)	SOKD (Student)
Top-1	26.69	30.25	29.34	29.30	29.38	30.04	28.83	28.38	29.13	25.85	<b>28.04</b>
Top-5	8.58	10.93	10.12	10.00	10.20	10.83	9.87	9.33	-	7.89	<b>8.92</b>

Table 3: Test error (%) of different offline methods on ImageNet. <sup>†</sup> SemCKD [8] only reports Top-1 error.

	Method	Vanilla	DML [21]	DCM [22]	SOKD
Teacher	ResNet50	23.87	24.23	<b>23.43</b>	23.49
Student	ResNet18	30.25	28.98	28.65	<b>28.44</b>

Table 4: Test error (%) of different online methods on ImageNet. We refer to the smaller network in online methods as student, otherwise as teacher.

training objectives can be modified as follows:

$$\begin{aligned}\mathcal{L}^{kbm} &= \alpha_1 \mathcal{L}_{ce}^{kbm} + \alpha_2 d(T_{kbm}(F_{kbm}), T_t(F_t)) + \alpha_3 d(T_{kbm}(F_{kbm}), T_s(F_s)), \\ \mathcal{L}^s &= \lambda_1 \mathcal{L}_{ce}^s + \lambda_2 d(T_s(F_s), T_{kbm}(F_{kbm})),\end{aligned}\quad (5)$$

where  $\alpha_1, \alpha_2, \alpha_3, \lambda_1, \lambda_2$  are balanced hyperparameters.  $F$  denotes the feature,  $T(\cdot)$  is feature transform function, and  $d(\cdot)$  is the distance function.

## 4 Experiments

### 4.1 Experimental Details

**Choices of Teacher and Student.** We consider six types of networks, including WideResNet [6], resnet [10], vgg [23], MobileNet [14], MobileNetV2 [24], and ShuffleNetV1 [11]. For convenience, we use WRN-d-w to denote the WideResNet with depth  $d$  and width factor  $w$  and MBNV2 ( $w$ ) to denote MobileNetV2 with a width multiplier of  $w$ . We use resnet $d$  and ResNet $d$  to represent CIFAR-style and ImageNet-style resnet with depth  $d$ , respectively. Besides, we use the following rules to construct teacher-student pairs: 1) the teacher and the student have the same architectural style (e.g., vgg13 and vgg8); 2) the teacher and the student are different architectures (e.g., WRN-28-4 and resnet56).

**Implementation Details.** We implement our methods on PyTorch [20]. We run compared methods based on two knowledge distillation benchmarks<sup>2</sup> and author-provided codes. 1) For CIFAR-100 dataset, we train all networks for 200 epochs. We use an SGD optimizer with a mini-batch size of 128, the momentum of 0.9, and the weight decay of  $5e-4$ . The learning rate is initialized by 0.1 (MobileNetV2 and ShuffleNetV1 by 0.05) and divided by 10 at both 100 and 150 epochs. 2) For ImageNet dataset, we follow the standard ImageNet parallel training practice on PyTorch. We set the batch size to 256 and train the student model for 100 epochs. The initial learning rate is 0.1 and decayed by 10 at 30, 60 and 90 epochs, respectively. We set  $\lambda_1 = \lambda_2 = \alpha_1 = \alpha_2 = \alpha_3 = 1$ ,  $\tau = 3$  for our SOKD on both CIFAR-100 and ImageNet. Code can be found at <https://github.com/swljq/Semi-Online-KD>.

<sup>2</sup><https://github.com/AberHu/Knowledge-Distillation-Zoo>; <https://github.com/HobbitLong/RepDistiller>

Method	Vanilla	AT [33]	CRD [26]	SSKD [52]	SemCKD [5]	KDCL [8]	DCM [54]	KD [42]	DML [42]	SOKD
CKA	0.7200	0.7287	0.7359	0.7773	0.7105	0.8815	0.8834	0.7966	0.8821	<b>0.8890</b>
IER (%)	23.99	23.81	21.31	19.87	24.18	20.39	16.89	22.77	17.15	<b>16.32</b>
MR (%)	<b>9.97</b>	10.03	11.09	11.59	10.13	11.04	12.91	10.18	13.24	11.61

Table 5: CKA similarity [47], imitation error rate (IER), and misleading rate (MR) on CIFAR-100 dataset. We take WRN-28-4 and WRN-16-2 as the teacher and the student, respectively.

## 4.2 Comparisons with Offline Methods

**Results on CIFAR-100.** In this experiment, we compare our proposed SOKD with KD [42], AT [33], SP [47], CC [20], VID [4], CRD [26], SSKD [52] and SemCKD [5] on eight teacher-student combinations. From Table 1, our SOKD outperforms all of compared methods by a large margin on all teacher-student pairs, which demonstrates the effectiveness and generalization of our method. It is worth noting that our SOKD not only achieves high performance in terms of the student’s accuracy, but also improves the performance of the teacher. Specifically, the teacher of our SOKD outperforms the vanilla pretrained teacher by 1.87% Top-1 accuracy, ranging from 1.54% to 2.32%. Similar to Vanilla KD method, our SOKD only leverages logit-based knowledge while showing superior performance than existing state-of-the-art feature-based methods. We show that the knowledge of logit still has huge room for improvement, which is usually underestimated by existing methods.

**Results on ImageNet.** In this experiment, we evaluate our proposed SOKD on ImageNet dataset. Constrained by limited computational resources, we only use ResNet34 and ResNet18 as teacher and student network, respectively. As shown in Table 3, our proposed SOKD outperforms other state-of-the-art methods. Specifically, our SOKD outperforms KD of 1.30% in terms of Top-1 accuracy. Besides, the teacher of SOKD achieves 0.84% higher Top-1 accuracy than vanilla teacher. These results show the generalization of our proposed SOKD in the large-scale dataset.

## 4.3 Comparisons with Online Methods

**Results on CIFAR-100.** In this experiment, we compare our SOKD with three state-of-the-art online KD methods (i.e., DML [42], KDCL [8], DCM [54]). We consider two training scenarios: 1) training two models with same backbone; 2) training two models with different backbones. We show the results in Table 2. From this table, we have the following observations: First, our SOKD gains the best student models. Second, it is interesting to find that SOKD further improves the accuracy of the teacher models, which never occurs in traditional offline KD methods. Specifically, our SOKD evenly achieves 1.96% higher Top-1 accuracy on four teacher-student pairs. Furthermore, although our main goal is to obtain a small and effective student model, we can further improve teacher by simply adjusting the structure of KBM, which will be discussed in section 4.4.

**Results on ImageNet.** In this experiment, we conduct experiments with two teacher-student pairs. We do not compare with KDCL [8] due to unavailable code on ImageNet. As shown in Table 4, our SOKD performs better than all compared methods. Specifically, for ResNet50-ResNet18 pair, SOKD has increased the accuracy of vanilla student and teacher by 1.81% and 0.38%, respectively. These results show effectiveness of our SOKD.



	Method	Vanilla	AT	DML-AT	SO-AT	SP	DML-SP	SO-SP	CRD	DML-CRD	SO-CRD
Teacher	vgg13	75.29	75.29	75.23	<b>75.88</b>	75.29	75.49	<b>76.68</b>	75.29	74.62	<b>76.41</b>
Student	vgg8	70.99	72.76	71.51	<b>73.04</b>	73.58	70.94	<b>74.29</b>	73.24	72.62	<b>73.74</b>
Teacher	WRN-28-4	78.91	78.91	78.78	<b>80.53</b>	78.91	79.53	<b>80.00</b>	78.91	79.03	<b>79.42</b>
Student	resnet56	73.05	73.25	73.26	<b>73.88</b>	74.28	73.15	<b>74.52</b>	74.19	73.85	<b>74.55</b>

Table 6: Test accuracy (%) of teacher and student on different training frameworks on CIFAR-100 dataset.

	Method	vanilla	KD [□]	classifier	classifier + 1 block	classifier + 2 blocks	classifier + 3 blocks	whole network
Teacher	WRN-28-4	78.91	78.91	79.35	80.47	80.64	<b>81.08</b>	80.96
Student	resnet56	73.05	74.21	74.88	<b>75.94</b>	75.90	75.87	75.34

Table 7: Effect of different KBMs on CIFAR-100 dataset. The last row “whole network” means the KBM is the same as the whole teacher network (It refers to WRN-28-4 here).

## 4.4 Further Experiments

**Imitation Ability of Student.** In this experiment, we try to analyze the degree to which the student imitates the teacher. From Table 5, we observe that 1) online methods (i.e., DML, KDCL, DCM) generally have learned similar representations to the teacher (higher CKA and lower IER) while suffering more serious misleading (higher MR). 2) our SOKD achieves much higher CKA and lower IER than all of compared methods, revealing that our SOKD has stronger ability to capture the behaviors of the teacher. Meanwhile, the supervision signals provided by the teacher in SOKD are much more accurate (lower MR) than that in DML. The main reason is that the KBM is able to obtain valuable information from the well-trained teacher.

**Extensions to Feature-Based Methods.** To investigate the scalability of our SOKD, we extend it to feature-based methods. Since feature information is related to network structure, we choose two different combinations of the teacher and the student in this experiment. Meanwhile, we also try to apply the training framework of DML into different methods. For simplicity, we use “DML/SO-M” to denote training M method with DML/SOKD framework. From Table 6, the performance of feature-based methods significantly decreases when applied with DML framework on most cases but our SOKD performs better than baselines. One possible reason is that in the early training of DML the representations might be meaningless, which would mislead the optimization directions of the teacher and the student. Instead, the training of SOKD is supervised by a well-trained teacher, which has learned powerful representations.

**Effect of Different KBMs.** In this experiment, we investigate the effect of different KBMs. We design the KBM with last classifier layers and last n feature extractor blocks of the teacher. As shown in Table 7, all KBMs achieve better performance than KD, which demonstrates the effectiveness and flexibility of our SOKD. The KBM can be changed according to different demands while maintaining the performance. However, the KBM with only classifier performs much worse than those with several extractor blocks, which means that it is necessary to adjust the parameters of intermediate layers such that the student is easier to learn similar representations from the teacher. However, taking the whole WRN-28-4 as KBM does not lead to further improvement, which shows the necessity of sharing teacher’s low-level layers. For the purpose of reducing computational cost, in this paper, we set the structure of KBM to be the same as the last feature extractor block and last classifier in most cases.

**Effect of Different Losses for KBM.** We investigate the effect of different losses for

$\mathcal{L}_{ce}^{kbm}$	$KL(p^{kbm}, p^t)$	$KL(p^{kbm}, p^s)$	WRN-28-4 (Teacher)	resnet56 (Student)
✓			78.02	75.45
✓	✓		79.24	74.91
✓		✓	79.34	74.87
	✓	✓	80.03	75.56
✓	✓	✓	<b>80.47</b>	<b>75.94</b>

Table 8: Effect of different losses for KBM on CIFAR-100 dataset.

KBM. The results are shown in table 8. When only equipped with  $\mathcal{L}_{ce}^{kbm}$ , the performance of the student is improved. The main reason is that KBM and student are trained simultaneously, the student is easier to imitate the KBM than the well-trained teacher. However, only adding the supervision of the teacher (i.e., w/o  $KL(p^{kbm}, p^s)$ ) or only obtaining the feedback from the student (i.e., w/o  $KL(p^{kbm}, p^t)$ ) both lead to performance degradation, which shows the indispensability of reliable signals and real-time feedback. Equipped with all losses, both teacher and student achieve highest accuracy, which shows the effectiveness of proposed losses for the KBM.

## 5 Conclusion

In this paper, we have discovered that KD and DML are complementary to each other on knowledge distillation task. Inspired by this observation, we have proposed a Semi-Online Knowledge Distillation (SOKD) method that integrates KD with DML. Our SOKD have effectively leveraged the well-trained teacher’s supervision signals and alleviated the imitation difficulty of the student. Moreover, we have extended the proposed framework to feature-based distillation methods, which also achieved competitive performance. Extensive experiments on two benchmarks have shown that our method consistently outperforms state-of-the-art methods on knowledge distillation tasks with different teacher-student pairs.

## Acknowledgement

This work was supported by the Key Realm R&D Program of Guangzhou (202007030007), China Scholarship Council (CSC).

## References

- [1] Sungsoo Ahn, Shell Xu Hu, Andreas Damianou, Neil D Lawrence, and Zhenwen Dai. Variational information distillation for knowledge transfer. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 9163–9171, 2019.
- [2] Saad Albawi, Tareq Abed Mohammed, and Saad Al-Zawi. Understanding of a convolutional neural network. In *2017 International Conference on Engineering and Technology (ICET)*, pages 1–6. IEEE, 2017.
- [3] Yoshua Bengio, Pascal Lamblin, Dan Popovici, and Hugo Larochelle. Greedy layer-wise training of deep networks. *Adv. Neural Inform. Process. Syst.*, 19:153–160, 2006.
- [4] Defang Chen, Jian-Ping Mei, Can Wang, Yan Feng, and Chun Chen. Online knowledge distillation with diverse peers. In *AAAI*, pages 3430–3437, 2020.

- [5] Defang Chen, Jian-Ping Mei, Yuan Zhang, Can Wang, Zhe Wang, Yan Feng, and Chun Chen. Cross-layer distillation with semantic calibration. *AAAI*, 2021.
- [6] Ting-Wu Chin, Ruizhou Ding, Cha Zhang, and Diana Marculescu. Towards efficient model compression via learned global ranking. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1518–1528, 2020.
- [7] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial networks. *Adv. Neural Inform. Process. Syst.*, abs/1406.2661, 2014.
- [8] J. Gou, Baosheng Yu, S. Maybank, and Dacheng Tao. Knowledge distillation: A survey. *Int. J. Comput. Vis.*, 2020.
- [9] Qiushan Guo, Xinjiang Wang, Yichao Wu, Zhipeng Yu, Ding Liang, Xiaolin Hu, and Ping Luo. Online knowledge distillation via collaborative learning. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 11017–11026, 2020.
- [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 770–778, 2016.
- [11] Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks. In *Int. Conf. Comput. Vis.*, pages 1389–1397, 2017.
- [12] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [13] Yuenan Hou, Zheng Ma, Chunxiao Liu, and Chen Change Loy. Learning lightweight lane detection cnns by self attention distillation. In *Int. Conf. Comput. Vis.*, pages 1013–1021, 2019.
- [14] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- [15] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 4700–4708, 2017.
- [16] Xiao Jin, Baoyun Peng, Yichao Wu, Yu Liu, Jiaheng Liu, Ding Liang, Junjie Yan, and Xiaolin Hu. Knowledge distillation via route constrained optimization. In *Int. Conf. Comput. Vis.*, pages 1345–1354, 2019.
- [17] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. *International Conference on Machine Learning*, 2019.
- [18] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Int. Conf. Comput. Vis.*, pages 2980–2988, 2017.

- [19] Yanxia Liu, Anni Chen, Hongyu Shi, Sijuan Huang, Wanjia Zheng, Zhiqiang Liu, Qin Zhang, and Xin Yang. Ct synthesis from mri using multi-cycle gan for head-and-neck radiation therapy. *Computerized medical imaging and graphics : the official journal of the Computerized Medical Imaging Society*, 91:101953, 2021.
- [20] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Adv. Neural Inform. Process. Syst.*, pages 8026–8037, 2019.
- [21] Baoyun Peng, Xiao Jin, Jiaheng Liu, Dongsheng Li, Yichao Wu, Yu Liu, Shunfeng Zhou, and Zhaoning Zhang. Correlation congruence for knowledge distillation. In *Int. Conf. Comput. Vis.*, pages 5007–5016, 2019.
- [22] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Adv. Neural Inform. Process. Syst.*, pages 91–99, 2015.
- [23] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *Int. Conf. Learn. Represent.*, 2015.
- [24] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 4510–4520, 2018.
- [25] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *Int. Conf. Learn. Represent.*, 2015.
- [26] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. *Int. Conf. Learn. Represent.*, 2020.
- [27] Frederick Tung and Greg Mori. Similarity-preserving knowledge distillation. In *Int. Conf. Comput. Vis.*, pages 1365–1374, 2019.
- [28] Hui Wang, Hanbin Zhao, Xi Li, and Xu Tan. Progressive blockwise knowledge distillation for neural network acceleration. In *IJCAI*, 2018.
- [29] Kuan Wang, Zhijian Liu, Yujun Lin, Ji Lin, and Song Han. Haq: Hardware-aware automated quantization with mixed precision. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 8612–8620, 2019.
- [30] Jiaxiang Wu, Cong Leng, Yuhang Wang, Qinghao Hu, and Jian Cheng. Quantized convolutional neural networks for mobile devices. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 4820–4828, 2016.
- [31] Zheng Xie, Zhiquan Wen, Jing Liu, Zhiqiang Liu, Xixian Wu, and Minghui Tan. Deep transferring quantization. In *Eur. Conf. Comput. Vis.*, 2020.
- [32] Guodong Xu, Ziwei Liu, Xiaoxiao Li, and Chen Change Loy. Knowledge distillation meets self-supervision. In *Eur. Conf. Comput. Vis.*, 2020.

- [33] Haoming Xu, Runhao Zeng, Qingyao Wu, Mingkui Tan, and Chuang Gan. Cross-modal relation-aware networks for audio-visual event localization. In *ACM Int. Conf. Multimedia*, pages 3893–3901, 2020.
- [34] Anbang Yao and Dawei Sun. Knowledge transfer via dense cross-layer mutual-distillation. In *Eur. Conf. Comput. Vis.*, pages 294–311. Springer, 2020.
- [35] Li Yuan, Francis EH Tay, Guilin Li, Tao Wang, and Jiashi Feng. Revisiting knowledge distillation via label smoothing regularization. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3903–3911, 2020.
- [36] Sukmin Yun, Jongjin Park, Kimin Lee, and Jinwoo Shin. Regularizing class-wise predictions via self-knowledge distillation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 13876–13885, 2020.
- [37] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *Brit. Mach. Vis. Conf.*, 2016.
- [38] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. *Int. Conf. Learn. Represent.*, abs/1612.03928, 2017.
- [39] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *Eur. Conf. Comput. Vis.*, pages 818–833. Springer, 2014.
- [40] Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be your own teacher: Improve the performance of convolutional neural networks via self distillation. In *Int. Conf. Comput. Vis.*, pages 3713–3722, 2019.
- [41] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 6848–6856, 2018.
- [42] Y. Zhang, T. Xiang, Timothy M. Hospedales, and H. Lu. Deep mutual learning. *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 4320–4328, 2018.
- [43] Xiatian Zhu, Shaogang Gong, et al. Knowledge distillation by on-the-fly native ensemble. In *Adv. Neural Inform. Process. Syst.*, pages 7517–7527, 2018.
- [44] Zhenzhou Zhuang, Zonghao Liu, Kin-Man Lam, Shuangping Huang, and Gang Dai. A new semi-automatic annotation model via semantic boundary estimation for scene text detection. In *ICDAR*, 2021.
- [45] Zhuangwei Zhuang, Mingkui Tan, Bohan Zhuang, Jing Liu, Yong Guo, Qingyao Wu, Junzhou Huang, and Jinhui Zhu. Discrimination-aware channel pruning for deep neural networks. In *Adv. Neural Inform. Process. Syst.*, pages 875–886, 2018.