# Multi-scale Feature Extraction and Fusion for Online Knowledge Distillation⋆

Panpan Zou[1][0000−0002−7040−6242], Yinglei Teng[1,2][0000−0002−7170−4764], and Tao Niu[1][0000−0001−6149−2908]

[1] Beijing University of Posts and Telecommunications, Beijing, China
{zoupanpan,lilytengtt,tasakim}@bupt.edu.cn
[2] Beijing Key Laboratory of Space-ground Interconnection and Convergence

**Abstract.** Online knowledge distillation conducts knowledge transfer among all student models to alleviate the reliance on pre-trained models. However, existing online methods rely heavily on the prediction distributions and neglect the further exploration of the representational knowledge. In this paper, we propose a novel Multi-scale Feature Extraction and Fusion method (MFEF) for online knowledge distillation, which comprises three key components: Multi-scale Feature Extraction, Dual-attention and Feature Fusion, towards generating more informative feature maps for distillation. The multi-scale feature extraction exploiting divide-and-concatenate in channel dimension is proposed to improve the multi-scale representation ability of feature maps. To obtain more accurate information, we design a dual-attention to strengthen the important channel and spatial regions adaptively. Moreover, we aggregate and fuse the former processed feature maps via feature fusion to assist the training of student models. Extensive experiments on CIFAR-10, CIFAR-100, and CINIC-10 show that MFEF transfers more beneficial representational knowledge for distillation and outperforms alternative methods among various network architectures.

**Keywords:** Knowledge distillation · Multi-scale · Feature fusion.

## 1 Introduction

Driven by the advances in algorithms, computing power, and big data, deep learning has achieved remarkable breakthroughs in various vision tasks [1,2,3]. Increasing the network depth or width is often the key point to further improve the performance of deep neural networks. However, these models with millions of parameters demand high computational costs and huge storage requirements, making it challenging to deploy them in resource-limited or low latency scenarios. For instance, mobile phones and Internet of Things (IoT) devices. To address this problem, extensive research has been carried out to develop a lightweight model while simultaneously keeping negligible model accuracy degradation in performance. These efforts can typically

---

be classified into network pruning, parameter quantization, low-rank approximation, and knowledge distillation.

Knowledge distillation (KD) has been demonstrated as an effective technique for model compression. The vanilla KD [4] method adopts a two-stage training strategy, where knowledge is transferred from the pre-trained high-capacity teacher model to a compact student model via aligning prediction distributions or feature representations [5], also known as the offline distillation. Drawbacks of these methods include the fact that the high-capacity teacher is not always available, even if they are, higher computational cost and training time of the complex teacher also cannot be avoided. In addition, KD suffers from model capacity gap when the size difference is large between the student and teacher model [6].

Online knowledge distillation (OKD) [7,8,9,10] has been developed to alleviate the above issue. This paradigm is more attractive for the reason that instead of using a pre-trained high-performance teacher, it breaks the presupposed specific strong-weak relationship and simplifies the training process to an end-to-end one-stage fashion. All models are trained simultaneously by learning from each other across the training process. In the other words, knowledge is distilled and shared among all networks. Compared to the offline KD, the online KD achieves superior performance while keeping a more straightforward structure. However, popular methods concentrate on transferring logit information as soft targets in common. Although the soft targets carry richer information than one-hot labels, it is relatively unitary to make use of only the logit. Since feature maps can provide rich information about the perception, channel and spatial correlations, simply aligning or fusing cannot take full advantage of the meaningful feature representation.

In this paper, to alleviate the aforementioned limitation, we propose a novel Multi-scale Feature Extraction and Fusion method (MFEF) for online knowledge distillation, including three key components, i.e., multi-scale feature extraction, dual-attention, and feature fusion. In order to obtain more beneficial representational knowledge for distillation, we first get multi-scale features which can focus on both local details and global regions by multiple divide and concatenate operations. Then, students are guided to learn more accurate features by introducing dual-attention which boosts the representation power of important channel and spatial regions while suppressing unnecessary regions. Finally, we utilize feature fusion to integrate the acquired feature maps and feed them into a fusion classifier to assist the learning of student models.

To summarize, the main contributions of this paper are:

- We propose a novel Multi-scale Feature Extraction and Fusion method (MFEF) for online knowledge distillation, which integrates the feature representation with soft targets for distillation.

- We first introduce multi-scale feature extraction to improve the multi-scale representation ability of the features and provide richer information apart from simply alignment. Then the dual-attention is proposed to generate more accurate features. Furthermore, we use feature fusion to fuse the enhanced knowledge, which can improve the generalization ability for distillation.

– Extensive experiments on CIFAR-10/100 [11] and CINIC-10 [12] verify that the proposed MFEF can effectively enhance the multi-scale representation power of features and generate more informative knowledge for distillation.

## 2   Related Work

Many efforts have been conducted with regard to knowledge distillation and vision tasks. In this section, we will give a comprehensive description of the related literature.
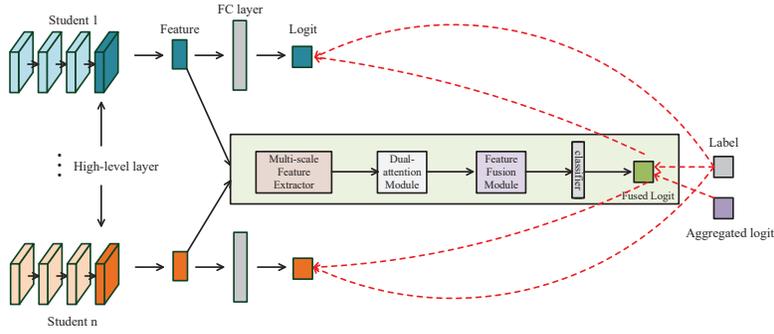
### 2.1   Traditional Knowledge Distillation

The idea of transferring the knowledge from a cumbersome model to a smaller model without a significant drop in accuracy is derived from [13]. Traditional KD works in a two-stage fashion which needs a pre-trained teacher. [5] first introduces intermediate features from hidden layers, the main idea is to match the feature activations of the student and teacher model. [14] combines attention with distillation to further exploit more accurate information. [15] explores the relationships between layers by mimicking the teacher's flow matrices using the inner product. In [16], the adversarial training scheme is utilized to enable the student and teacher networks to learn the true data distribution. [6] introduces a teacher assistant to mitigate the capacity gap between the teacher model and student model. In [17], it proposes to use the activation boundaries formed by hidden neurons for distillation.

### 2.2   Online Knowledge Distillation

Online knowledge distillation has emerged to further improve the performance of the student model and eliminate the dependency on high-capacity teacher models which are time-consuming and costly. In this paradigm, student models learn mutually by sharing the predictions throughout the training process. [7] is a representative method in which multiple networks work in a collaborative way. Each network imitates the peer network's class probabilities using Kullback-Leibler divergence. [9] further extends DML to construct an ensemble logit as the teacher by averaging a group of students' predictions to improve generalization ability. A fusion module is introduced to train a fusion classifier to guide the training of sub-networks in [10]. [18] adds a gate module to generate the importance score for each branch and build a stronger teacher. [19] proposes two-level distillation between multiple auxiliary peers and a group leader to enhance diversity among student models. In terms of architecture designing, [20] forms the student model via replacing the standard convolution with cheap convolution operations. Student and teacher models share the same networks in [21], where knowledge is distilled within the network itself and knowledge from the deeper portions of the network is distilled into shallow ones.

### 2.3   Multi-scale Feature

Multi-scale feature representations are of critical importance to many vision tasks. Some concurrent works focus on promoting the capability of models by utilizing

**Fig. 1.** An overview of Multi-scale Feature Extraction and Fusion (MFEF) for Online knowledge distillation. The output of high-level layer is sent to three key components (i) Multi-scale Feature Extraction: Enhance the multi-scale representation ability of feature maps. (ii) Dual-attention: Use channel and spatial attention to strengthen informative regions. (iii) Feature Fusion: Integrate knowledge among stuent models and futher improve the generalization ability.

multi-scale features. [22] constructs hierarchical residual-like connections within a residual block to represent multi-scale features at a granular level. [23] uses pyramidal convolution including four levels of different kernel sizes to generate multi-scale features. Similarly, [24] integrates information at different scales via pyramidal convolution structure for the channel-wise feature maps. A flexible and efficient hierarchical-split block is used in [25] to capture multi-scale features. [26] adopts atrous spatial pyramid pooling to probes convolutional features on multiple scales for semantic image segmentation.
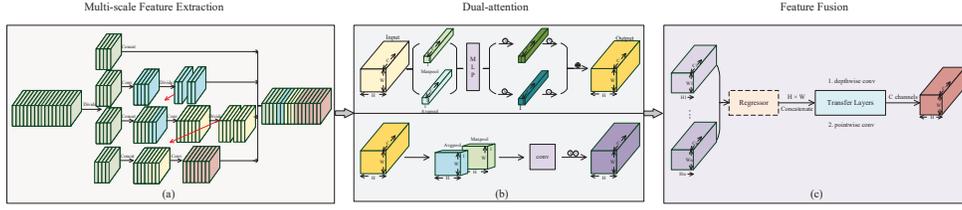
## 3   Proposed Method

In this section, we describe the framework and loss function in detail. An overview of MFEF is illustrated in Fig. 1. Different from the existing KD methods, MFEF digs deeper into the information provided by feature maps including multi-scale representation ability and the channel and spatial attention.

### 3.1   Problem Definition

The key idea of knowledge distillation is that soft targets contain the dark knowledge which can be used as a supervisor to transfer knowledge to the student model. Given a labeled dataset $D\{x_i, y_i\}_{i=1}^N$, with $N$ samples, $x_i$ is the $i$th input sample and $y_i \in \{1, 2, ..., M\}$ is the corresponding ground-truth label. $M$ is the total number of classes in the dataset. Consider $n$ student models $\{S_j\}_{j=1}^n$, the logit produced by the last fully connected layer of the student $S_j$ is denoted as $z_j = \{z_j^1, z_j^2, ..., z_j^M\}$. Then the probability of the $j$th student for the sample $x_i$ over the $m$th class $p_j^m(x_i)$ can be estimated by a softmax function,

$$p_j^m(x_i) = \frac{\exp(z_j^m/T)}{\sum_{m=1}^M \exp(z_j^m/T)}, \tag{1}$$

**Fig. 2.** The structure of key components: (a). Multi-scale feature extraction. (b) Dual-attention. (c) Feature Fusion

where $T$ is the temperature which produces a more softened probability distribution as it increases. Specifically, when $T = 1$, it is defined as the original softmax output, we consider writing it as $p_j^m(x_i)$; otherwise it is rewritten as $\tilde{p}_j^m(x_i)$. For multi-class classification, the objective is to minimize the cross-entropy loss between the softmax outputs and the ground-truth labels,

$$L_j^{CE} = -\sum_{i=1}^{N}\sum_{m=1}^{M} l_i \log(p_j^m(x_i)), \tag{2}$$

where $l_i = 1$ if $y_i = m$, and 0 otherwise. Knowledge transfer is facilitated by matching the softened probability of the student model $\tilde{p}_j^m(x_i)$ and the teacher model $\tilde{p}_t^m(x_i)$. We introduce the distillation loss of $j - th$ student model in the form of Kullback-Leibler Divergence

$$L_j^D = \sum_{i=1}^{N}\sum_{j=1}^{M} \tilde{p}_t^m(x_i) \log \frac{\tilde{p}_t^m(x_i)}{\tilde{p}_j^m(x_i)}. \tag{3}$$

### 3.2   MFEF Framework

From a global perspective, the main idea of MFEF is to enhance the multi-scale representation power of feature maps and generate more informative knowledge for distillation. The details of each key component are explained in the following.

**Multi-scale Feature Extraction.** Aligning the soft targets of teacher and student models enhances the model generalization, but it ignores the feature maps which contain rich information. In addition to the soft targets, inspired by [25], we introduce multi-scale feature extraction to generate multi-scale features which are of significant importance for vision tasks. As shown in Fig. 2 (a), the extraction includes multiple divide and concatenate operations in the channel dimension to enhance the information flow between different groups. We use the feature maps of the last layer as the input for the reason it has high-level semantic information which is richer and specific. For the convenience of notation, we name the feature map of the $j$th student model as $F_j$ and the extraction as $E$. $D$ and $C$ represent the divide and concatenate operations, respectively. First, $F_j$ is divided into $p$ groups $\{F_{j1}, F_{j2}, ..., F_{jp}\}$. The first group $F_{j1}$ is output straightforward and the second part is sent to a convolution operation and then is divided into two sub-groups $D_{21}$ and $D_{22}$. One of them is exported

to the output and the other is concatenated with the next part. The rest other than the last group follows the concatenate-convolution-export-divide procedure. The last part does not need the divide operation. We define the output as

$$E(F_j) = C(F_{j1}, D_{22}, D_{32}, ..., Conv(C(F_{jp}, D_{p-1,2}))). \qquad (4)$$

The multi-scale feature extraction can generate feature maps that contain multiple scales of receptive fields. The more features are concatenated, the larger the receptive field is. Larger receptive fields can capture global information while the smaller ones can focus on details. Such a combination can generate more meaningful feature maps to improve the performance of distillation.

**Dual-attention.** After the extraction, we utilize dual-attention to dig deeper into the feature maps (see Fig. 2 (b)). Channel and spatial attention focus on "what" and "where" are important, and we apply them in a sequential manner. We denote the multi-scale feature map $E(F_j) \in \mathbb{R}^{C \times H \times W}$ as the input, where $C$, $H$, $W$ represent its channel numbers, height, and width, respectively. Average-pooling and max-pooling are used in combination to obtain finer attention.

For channel attention, we denote $a_c, m_c \in \mathbb{R}^{C \times 1 \times 1}$ as the vectors after average-pooling and max-pooling. The weight $w_c \in \mathbb{R}^{C \times 1 \times 1}$ of channel is

$$w_c = \sigma(W(a_c))) + (W(m_c))), \qquad (5)$$

where the symbol $\sigma$ denotes the Sigmoid function, $W$ is the weight of a multi-layer perceptron (MLP). The channel attention output $AT_j^c$ is

$$AT_j^c = w_c \otimes F_M, \qquad (6)$$

where $\otimes$ refers to element-wise multiplication. Similarly, we denote the average-pooling and max-pooling vector $a_s, m_s \in \mathbb{R}^{1 \times H \times W}$, $w_s \in \mathbb{R}^{1 \times H \times W}$ is

$$w_s = \sigma(conv(a_s; m_s)), \qquad (7)$$

where *conv* represents a convolution operation. The output $AT_j^s$ is

$$AT_j^s = w_s \otimes AT_j^c, \qquad (8)$$

Dual-attention can strengthen the informative channel and spatial regions while suppressing the less important ones and thus generate more informative outputs which can focus on useful regions within a context adaptively.

**Feature Fusion.** We propose feature fusion to aggregate and maximize the usage of the student models' information. The structure of it is illustrated in Fig. 2 (c). Specifically, we first concatenate the meaningful feature maps of students that have been processed previously, i.e., $\{AT_1^s, AT_2^s, ..., AT_j^s\}$. If the resolutions of the feature maps are different, we apply a convolutional regressor to make them identical. Then we concatenate them and sent the results to the transfer layers which consist of a sequence of depthwise and pointwise convolution operations. Finally, we fuse the student models' feature information and feed it into a fusion classifier which is supervised by ground truth labels.

### 3.3   Loss Function

The cross-entropy loss of the $j$th student and the fused classifier is $L_j^{CE}$ and $L_f^{CE}$, respectively, as described in Eq. (2). We further define the aggregated logit of students as $z_a^m = \frac{1}{n}\sum_{j=1}^n z_j^m$ and probability as $p_a^m$. The fusion classifier is trained with KL divergence

$$L_a^D = L_a^{KL}(\tilde{p}_a^m, \tilde{p}_f^m), \tag{9}$$

This loss is used to transfer the knowledge of the student models to the fusion classifier. Then the fusion classifier facilitate the knowledge which contains informative feature representations transferring back to the student models via minimizing the distillation loss

$$L_f^D = \sum_{j=1}^n L_f^{KL}(\tilde{p}_f^m, \tilde{p}_j^m), \tag{10}$$

Finally, we derive the total training objective as

$$L_{total} = L_{CE} + T^2 L_D. \tag{11}$$

where $L_{CE}$ is the sum of cross-entropy of students and fused classifier. $L_D$ refers to the sum of $L_a^D$ and $L_f^D$. Because the gradients produced by the soft targets are scaled by $1/T^2$, thus $L_D$ is multiplied with $T^2$ to keep the contributions of $L_{CE}$ and $L_D$ roughly balanced.

## 4   Experiment

In this section, we conduct comprehensive experiments to evaluate the performance of MFEF on three datasets and various widely-used neural networks. We choose various related methods under different settings for comparison and show the results to demonstrate that MFEF generalizes well among different numbers and types of models. Finally, evaluation of each component are carried out.

### 4.1   Experiment Settings

**Datasets and Architecture.** We incorporate three image classification datasets in the following evaluations. (1) CIFAR-10 which contains 60000 colored natural images (50000 training samples and 10000 test samples) over 10 classes. (2) CIFAR-100 consists of 60000 images (50000 training samples and 10000 test samples) drawn from 100 classes. (3) CINIC-10 consists of images from both CIFAR and ImageNet. It is more challenging than CIFAR-10. It contains 90000 train samples and 90000 test samples. For CIFAR-10/100, there are seven popular networks used, namely ResNet-20, ResNet-32, ResNet-56, ResNet-110, WRN-16-2, WRN-40-2, and DenseNet-40-12. For CINIC-10, we use MobileNetV2 and ResNet-18 following the settings in [12].

**Table 1.** Comparisons with closely related methods on CIFAR 10 and CIFAR-100 with seven different networks. Top-1 error rates(%) are reported. Two same student models are used for each method. FFL-S and MFEF-S refer to the results of the student model, and FFL and MFEF refer the results of fused classifiers.
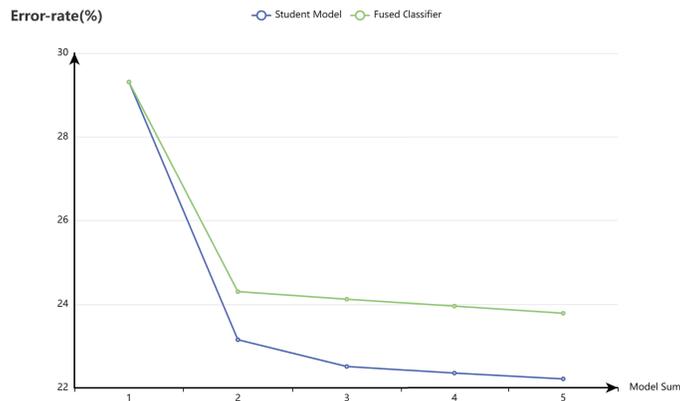
| Dataset | Network | Baseline | KD | DML | FFL-S | FFL | MFEF-S | MFEF |
|---------|---------|----------|-----|-----|-------|-----|--------|------|
| CIFAR-10 | ResNet-20 | 7.32 | 7.18 | 6.63 | 6.49 | 6.22 | **6.38** | **6.08** |
| | ResNet-32 | 6.77 | 6.69 | 6.52 | 6.06 | 5.78 | **5.59** | **5.41** |
| | ResNet-56 | 6.30 | 6.14 | 5.82 | 5.46 | 5.26 | **5.28** | **4.82** |
| | ResNet-110 | 5.64 | 5.47 | 5.21 | 5.18 | 4.83 | **4.81** | **4.52** |
| | WRN-16-2 | 6.78 | 6.40 | 5.49 | 6.09 | 5.97 | **5.33** | **4.99** |
| | WRN-40-2 | 5.34 | 5.24 | 4.72 | 4.75 | 4.60 | **4.51** | **4.02** |
| | DenseNet40-12 | 6.87 | 6.81 | 6.50 | 6.72 | 6.24 | **5.79** | **5.30** |
| CIFAR-100 | ResNet-20 | 31.08 | 29.94 | 29.61 | 28.56 | 26.87 | **28.46** | **26.30** |
| | ResNet-32 | 30.34 | 29.82 | 26.89 | 27.06 | 25.56 | **26.36** | **24.84** |
| | ResNet-56 | 29.31 | 28.61 | 25.51 | 24.85 | 23.53 | **24.22** | **23.15** |
| | ResNet-110 | 26.30 | 25.67 | 24.49 | 23.95 | 22.79 | **23.37** | **22.16** |
| | WRN-16-2 | 27.74 | 26.78 | 26.16 | 25.72 | 24.74 | **24.66** | **22.93** |
| | WRN-40-2 | 25.13 | 24.43 | 22.77 | 22.06 | 21.05 | **21.76** | **20.60** |
| | DenseNet40-12 | 28.97 | 28.74 | 26.94 | 27.21 | 24.76 | **26.81** | **24.27** |

**Settings.** We apply horizontal flips and random crop from an image padded by 4 pixels for data augmentation in training. We use SGD as the optimizer with Nesterov momentum 0.9, weight decay of 1e-4 for student models and 1e-5 for fusion and mini-batch size of 128. The models are trained for 300 epochs for all datasets. We set the initial learning rate to 0.1 and is multiplied by 0.1 at 150, 225 epochs. We set the temperature $T$ to 3 empirically and $\alpha = 80$ for ramp-up weighting. For the case of student models have same architecture, the low-level layers are shared following [19]. When output channels of the feature maps are different, the feature fusion is designed to match the smaller one. For fair comparison, we set the number of student models to two. The top-1 error rate (%) of the best student over 3 runs is reported.

### 4.2   Experiment Results

**Results on CIFAR-10/100.** As shown in Table 1, we evaluate the effectiveness of MFEF on CIFAR-10 and CIFAR-100 based on several popular networks. Since our goal is to distill more powerful feature representations for online distillation, we compare MFEF with the offline KD, logit-only online method DML, and fusion-only method FFL. For the offline KD, it employs a pre-trained ResNet-110 as the teacher model. For DML, we report the top-1 error rate of the best student. FFL-S and MFEF-S represent the results of the best student and FFL and MFEF indicate the results of the fused classifier.

The results clearly show the performance advantages of our MFEF. Specifically, MFEF improves by approximately 1% and 2% of the backbone networks. MFEF also achieves the best top-1 error rate compared with the closely related online distillation

**Fig. 3.** Evaluating the impact of expansion of student models on CIFAR-100 using ResNet-56.

methods. For instance, on CIFAR-10, MFEF-S achieves lower error rates than FFL-S by approximately 0.5%, 0.8%, and 1% on ResNet-32, WRN-16-2, and DensNet-40-12, respectively. MFEF improves FFL by about 1% on WRN-16-2 and DensNet-40-12; While on CIFAR-100, MFEF achieves 0.6%, 0.7%, and 1% increase on ResNet-56, ResNet -32, and WRN-16-2, respectively. MFEF is higher by about 1.8% on WRN-16-2 compared with FFL. These improvements attribute to the integration of the multi-scale feature extraction and the attention mechanism and the feature fusion of student models.

**Table 2.** Top-1 error rate (%) comparison with FFL on CINIC-10.

| Network | Baseline | FFL-S | FFL | MFEF-S | MFEF |
|---|---|---|---|---|---|
| MobileNetV2 | 18.07 | 17.85 | 16.10 | **17.56** | **15.66** |
| ResNet-18 | 13.94 | 13.33 | 12.67 | **13.22** | **12.39** |

**Results on CINIC-10.** In this section, we compare the top-1 error rate of MFEF with FFL based on MobileNetV2 and ResNet-18. As shown in Table 2, FFL and MFEF both reduces the error rate of the baseline and MFEF shows higher improvement of performance in both student models and fused classifier. In case of MobileNetV2, MFEF improves by around 0.5%, 0.3%, and 0.4% compared to the baseline and FFL for student model and fused classifier. Based on these experiments, we could confirm that thanks to the enhancement of the multi-scale representation power, higher-quality knowledge is transferred among all student models and consequently achieves a lower error rate than others.

**Expansion of Student Models.** The impact of increasing the number of student models is illustrated in Fig.3. We conduct experiments on ResNet-56. Not surpris-

**Table 3.** Top-1 error rate (%) comparisons with other online distillation methods for training three students model on CIFAR-100. ONE and ONE-E refer to the results of the student models and the gated ensemble teacher.

| Method | ResNet-32 | ResNet-56 |
|--------|-----------|-----------|
| ONE    | 26.64     | 24.63     |
| FFL-S  | 26.30     | 24.51     |
| MFEF-S | **26.04** | **24.12** |
| ONE-E  | 24.75     | 23.27     |
| FFL    | 24.31     | 23.20     |
| MFEF   | **24.03** | **22.51** |

**Table 4.** Top-1 error rate (%) comparisons with other online distillation methods for different architectures of student models on CIFAR-100.

|      | Net1 ResNet-32 | Net2 WRN-16-2 | Net1 ResNet-56 | Net2 WRN-40-2 |
|------|----------------|---------------|----------------|---------------|
| DML  | 28.31          | 26.45         | 26.75          | 23.33         |
| FFL  | 27.06          | 25.93         | 26.23          | 23.06         |
| MFEF | **26.38**      | **25.16**     | **25.7**       | **22.39**     |

ingly, the performance of both students and the fusion classifier improves as the number of student models increases. When the student models expanded to 3, MFEF still performs competitively against ONE and FFL, as shown in Table 3. We can see that MFEF-S achieves an approximately 0.3% and 0.6% performance improvement on ResNet-32 compared to FFL-S and ONE-S, respectively. The fusion classifier yields an about 0.7% and 0.8% improvement on ResNet-56 superior to ONE and FFL.

**Different Architecture.** To verify the generalization of MFEF on different model architectures, we conduct experiments on ResNet and WRN in Table 4. We set ResNet as Net1 and WRN as Net2. MFEF shows better performance than DML and FFL in both Net1 and Net2. An interesting observation is that when MFEF is applied, the smaller network (Net1) improves significantly compared to the larger one. For example, when compared with DML, MFEF is higher by about 2% and 1.3% on ResNet-32 and WRN-16-2. This is because MFEF can aggregate and fuse all networks' feature maps and transfer the informative knowledge of the larger network to the smaller one better.

**Table 5.** Evaluating the effectiveness of each component on CIFAR-100 using ResNet-110.

| Case | Component            | Student | Fused |
|------|----------------------|---------|-------|
| A    | Backbone             | 26.30   | -     |
| B    | Backbone+MSFE        | 24.35   | -     |
| C    | Backbone+OKD         | 24.79   | -     |
| D    | Backbone+MSFE+OKD    | 23.54   | 22.54 |
| E    | Backbone+MSFE+DA+OKD | 23.37   | 22.16 |

**Component Effectiveness Evaluation.** To further validate the benefit of each component, we conduct various ablation studies on CIFAR-100 on ResNet-110. Specifically, we perform experiments in five cases of ablations. As shown in Table 5, Case A refers to the model trained from scratch. Case B and C refer to the network where only the multi-scale feature extraction (MSFE) and OKD are included. And they improved by around 2% and 1.5% compared to the backbone. When both MSFE and OKD are applied in Case D, the student model achieves a higher accuracy by around 2.8% compared to Case A. When we get rid of MSFE from Case D (Case C), the performance decrease sharply by about 1.3%, which confirms the usefulness of the MSFE. Dual-attention (DA) is added in Case E based on Case D. This increases the performance by around 0.2% and 0.4% of the student models and the fused classifier, respectively, and it has more influence on the fused classifier. The improvements manifest that MSFE has a more significant impact on the model performance, which is mainly attributed to the enhancement of the multi-scale representation ability.

## 5   Conclusion

We present a novel Multi-scale Feature Extraction and Fusion method (MFEF) for online knowledge distillation . It integrates multi-scale extraction and attention mechanism into a unified feature fusion framework. Different from existing online knowledge distillation methods, we enhance the multi-scale representation ability of the feature maps and then fuse them from student models to assist the training process by transferring more informative knowledge. Extensive experiments based on three datasets show the superiority of our method compared to prior works. Results on various networks also demonstrate that the proposed method can be broadly applied to a variety of architectures from a very small scale to a large one.

## References

1. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. Advances in neural information processing systems **25** (2012)
2. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. Advances in neural information processing systems **28** (2015)
3. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3431–3440 (2015)
4. Hinton, G., Vinyals, O., Dean, J., et al.: Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531 **2**(7) (2015)
5. Romero, A., Ballas, N., Kahou, S.E., Chassang, A., Gatta, C., Bengio, Y.: Fitnets: Hints for thin deep nets. arXiv preprint arXiv:1412.6550 (2014)
6. Mirzadeh, S.I., Farajtabar, M., Li, A., Levine, N., Matsukawa, A., Ghasemzadeh, H.: Improved knowledge distillation via teacher assistant. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 5191–5198 (2020)
7. Zhang, Y., Xiang, T., Hospedales, T.M., Lu, H.: Deep mutual learning. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4320–4328 (2018)

8. Zhou, G., Fan, Y., Cui, R., Bian, W., Zhu, X., Gai, K.: Rocket launching: A universal and efficient framework for training well-performing light net. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 32 (2018)

9. Guo, Q., Wang, X., Wu, Y., Yu, Z., Liang, D., Hu, X., Luo, P.: Online knowledge distillation via collaborative learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 11020–11029 (2020)

10. Kim, J., Hyun, M., Chung, I., Kwak, N.: Feature fusion for online mutual knowledge distillation. In: 2020 25th International Conference on Pattern Recognition (ICPR). pp. 4619–4625. IEEE (2021)

11. Krizhevsky, A., Hinton, G., et al.: Learning multiple layers of features from tiny images (2009)

12. Darlow, L.N., Crowley, E.J., Antoniou, A., Storkey, A.J.: Cinic-10 is not imagenet or cifar-10. arXiv preprint arXiv:1810.03505 (2018)

13. Buciluǎ, C., Caruana, R., Niculescu-Mizil, A.: Model compression. In: Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining. pp. 535–541 (2006)

14. Zhou, Z., Zhuge, C., Guan, X., Liu, W.: Channel distillation: Channel-wise attention for knowledge distillation. arXiv preprint arXiv:2006.01683 (2020)

15. Yim, J., Joo, D., Bae, J., Kim, J.: A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4133–4141 (2017)

16. Wang, X., Zhang, R., Sun, Y., Qi, J.: Kdgan: Knowledge distillation with generative adversarial networks. Advances in Neural Information Processing Systems **31** (2018)

17. Heo, B., Lee, M., Yun, S., Choi, J.Y.: Knowledge transfer via distillation of activation boundaries formed by hidden neurons. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 33, pp. 3779–3787 (2019)

18. Zhu, X., Gong, S., et al.: Knowledge distillation by on-the-fly native ensemble. Advances in neural information processing systems **31** (2018)

19. Chen, D., Mei, J.P., Wang, C., Feng, Y., Chen, C.: Online knowledge distillation with diverse peers. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 3430–3437 (2020)

20. Xie, J., Lin, S., Zhang, Y., Luo, L.: Training convolutional neural networks with cheap convolutions and online distillation. arXiv preprint arXiv:1909.13063 (2019)

21. Zhang, L., Song, J., Gao, A., Chen, J., Bao, C., Ma, K.: Be your own teacher: Improve the performance of convolutional neural networks via self distillation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 3713–3722 (2019)

22. Gao, S.H., Cheng, M.M., Zhao, K., Zhang, X.Y., Yang, M.H., Torr, P.: Res2net: A new multi-scale backbone architecture. IEEE transactions on pattern analysis and machine intelligence **43**(2), 652–662 (2019)

23. Duta, I.C., Liu, L., Zhu, F., Shao, L.: Pyramidal convolution: Rethinking convolutional neural networks for visual recognition. arXiv preprint arXiv:2006.11538 (2020)

24. Zhang, H., Zu, K., Lu, J., Zou, Y., Meng, D.: Epsanet: An efficient pyramid split attention block on convolutional neural network. arxiv 2021. arXiv preprint arXiv:2105.14447

25. Yuan, P., Lin, S., Cui, C., Du, Y., Guo, R., He, D., Ding, E., Han, S.: Hsresnet: Hierarchical-split block on convolutional neural network. arXiv preprint arXiv:2010.07621 (2020)

26. Redmon, J., Farhadi, A.: Yolov3: An incremental improvement. arXiv preprint arXiv:1804.02767 (2018)