

CDFKD-MFS: Collaborative Data-free Knowledge Distillation via Multi-level Feature Sharing

Zhiwei Hao, Yong Luo, *Member, IEEE*, Zhi Wang *Member, IEEE*, Han Hu, *Member, IEEE*,
and Jianping An, *Member, IEEE*,

Abstract—Recently, the compression and deployment of powerful deep neural networks (DNNs) on resource-limited edge devices to provide intelligent services have become attractive tasks. Although knowledge distillation (KD) is a feasible solution for compression, its requirement on the original dataset raises privacy concerns. In addition, it is common to integrate multiple pretrained models to achieve satisfactory performance. How to compress multiple models into a tiny model is challenging, especially when the original data are unavailable. To tackle this challenge, we propose a framework termed collaborative data-free knowledge distillation via multi-level feature sharing (CDFKD-MFS), which consists of a multi-header student module, an asymmetric adversarial data-free KD module, and an attention-based aggregation module. In this framework, the student model equipped with a multi-level feature-sharing structure learns from multiple teacher models and is trained together with a generator in an asymmetric adversarial manner. When some real samples are available, the attention module adaptively aggregates predictions of the student headers, which can further improve performance. We conduct extensive experiments on three popular computer visual datasets. In particular, compared with the most competitive alternative, the accuracy of the proposed framework is 1.18% higher on the CIFAR-100 dataset, 1.67% higher on the Caltech-101 dataset, and 2.99% higher on the mini-ImageNet dataset.

Index Terms—Model Compression, Knowledge Distillation, Data-free Distillation, Multi-teacher Distillation, Attention

I. INTRODUCTION

DURING the past few years, deep neural networks (DNNs) have become one of the most influential models in machine learning, and with the developments of the 5G wireless network and artificial intelligence of things (AIoT), equipping resource-constrained edge devices with DNN-based intelligent applications has become appealing. Unfortunately, modern DNNs are heavily overparameterized, requiring tremendous storage and computing resources for inference. For example, ResNet-50 [1], a typical DNN architecture for computer vision, occupies more than 100 megabytes of ROM for storage and requires more than 4 giga FLOPs for inference. Such resource greedy DNNs make deployment challenging. Knowledge distillation (KD) [2] provides a feasible solution

by compressing a large pretrained teacher model (teacher) into a tiny student model (student). However, the dataset for training the large model, which is indispensable in KD, is sometimes unavailable due to the high transmitting overhead or privacy concerns. Although some data-free KD approaches have been proposed [3], [4], they are not suitable for multi-model scenarios, where multiple models are integrated in an ensemble to achieve better performance. How to effectively compress multiple models by KD in a data-free manner is still an open problem.

There are some obstacles to achieving a multi-teacher data-free KD. First, the student may be confused from the knowledge provided by different teachers. Since the benefit of an ensemble model usually comes from its diverse feature representations learned by multiple models [5], a student with a linear-structure that is directly trained using these models may not be able to sufficiently preserve diversity, which results in performance degradation. Second, conducting KD requires careful design. Directly regarding the ensemble as a single model to apply single-teacher KD is convenient, but this method may ignore the diversity of different teachers and lead to a suboptimal result. Finally, how to effectively leverage available samples from the original dataset for data-free KD for better performance is challenging.

Some works designed multi-branch student architectures to avoid knowledge confusion when the original dataset is available. For example, Tran *et al.* [6] proposed a tree-like student architecture for multi-teacher KD, and Liu *et al.* [7] designed a multi-task attention network to learn multiple tasks simultaneously. A few works studied how to achieve KD with multiple teachers when the original data are unavailable. Ye *et al.* [8] proposed a group-stack dual-generative adversarial network (GAN) architecture. However, their work mainly focused on the knowledge amalgamation (KA) problem. Their model aimed to amalgamate functions of multiple models in one model and little attention was paid to the size of the target model. Specifically, a linear network was trained as a backbone and the last several layers of teachers was grafted the onto the backbone. This approach achieved KA but did not compress the teacher branches. Moreover, the backbone of the model is a concatenation of multiple separately trained parts, which may not reach the global optima. For the second challenge, generator-based approaches [3], [9] provide a favorable solution for the single-teacher scenario. Furthermore, in data-free quantization, where the teacher and the student have the same architecture, intermediate feature supervision has been introduced [10]. However, it is still

Corresponding author: Han Hu.

Zhiwei Hao, Han Hu, and Jianping An are with the School of Information and Electronics, Beijing Institute of Technology, Beijing 100081, China (e-mail: haozhw@bit.edu.cn, hhu@bit.edu.cn, an@bit.edu.cn).

Yong Luo is with the School of Computer Science, Wuhan University, Wuhan 430072, China (e-mail: luoyong@whu.edu.cn).

Zhi Wang is with Shenzhen International Graduate School, Tsinghua University, Shenzhen 518055, China (e-mail: wangzhi@sz.tsinghua.edu.cn).

The code is available at <https://github.com/Hao840/CDFKD-MFS>

desirable to design an approach to conduct data-free KD on multiple teachers and a student with a different architecture. Finally, using real samples in data-free KD has not been studied. Because a potential domain gap exists between real and generated samples, directly combining them to train a student is inappropriate.

To tackle these issues, we propose a collaborative data-free knowledge distillation framework with a multi-level feature-sharing (CDFKD-MFS) student architecture. This framework is an extension of our previous conference work [11]. Our new CDFKD-MFS framework consists of three parts: a multi-header student network module, an asymmetric adversarial data-free KD module, and an attention-based prediction aggregation module. First, the student is composed of a backbone network and multiple lightweight header networks and takes the multi-level feature of the backbone as input. This architecture alleviates the knowledge confusion problem in multi-teacher KD by decoupling parameters for learning shared knowledge and teacher-specific knowledge. Then, the student and a data generator are adversarially trained based on the distillation loss. In addition, the generator is guided by statistics in the batch normalization layers of teachers, and the student is trained using an ensemble distillation loss and a feature distillation loss. Finally, for cases when some real samples are available, an attention module is designed to aggregate predictions of multiple headers in the student. We evaluate the proposed framework on three datasets with different sample sizes: CIFAR-100 [12], Caltech-101 [13], and mini-ImageNet [14]. The results show that our method achieves at least 1.18%, 1.67%, and 2.99% performance improvements compared to the next best model on these datasets, respectively. We summarize our contributions as follows:

- We propose a data-free KD framework for multi-model compression, where multiple pretrained teachers collaborate to obtain a better student.
- We design a multi-header student equipped with a multi-level feature-sharing architecture, where its backbone learns shared knowledge, and each header learns from a specific teacher.
- We design an asymmetric adversarial data-free KD algorithm, where a data generator and the student are trained in an asymmetric adversarial way.
- We develop an attention-based prediction aggregation module, which can further improve performance when a few samples from the original dataset are available.

The rest of this paper is organized as follows. Section II introduces existing works related to KD, branchy networks, and attention. Section III provides details of the proposed framework. We evaluate the framework in Section IV and finally summarize our work in Section V.

II. RELATED WORKS

A. Data-required KD

Knowledge distillation, which is a practical approach for compressing a large pretrained teacher model into a tiny student model, was first proposed by Hinton *et al.* [2]. The original KD method takes predictions of the teacher model

as soft labels to train the tiny student model from a sketch. Based on this framework, some works simulate soft labels by regularizing [15] or label smoothing [16], [17]. In addition to the output layer, there are also some approaches for transferring the knowledge in intermediate layers of the teacher. Based on this consideration, FitNet [18] and some follow-up works [19], [20] adopted feature matching losses to utilize intermediate knowledge. Moreover, intermediate-layer supervision can also be achieved through attention map matching [21], feature distribution matching [22], and cross-layer knowledge transferring [23].

Multiple pretrained teacher models can provide diverse knowledge and improve the performance of the student [5]. Hinton *et al.* [2] proposed using the average of teacher prediction as the supervising signal. You *et al.* [24] and Park *et al.* [25] introduced feature supervision in multi-teacher KD. To better transfer knowledge from multiple teachers, some works added additional branches to the student [6], [26]. However, these branches were directly grafted onto the student without careful consideration, which may not be sufficient for preserving the diverse knowledge of teachers. Hence, multi-teacher KD requires a student to learn diverse knowledge from multiple teachers while remaining lightweight. To tackle this problem, we design a multi-header student with a multi-level feature-sharing structure to learn from multiple teachers.

B. Data-free KD

Data-free KD aims to achieve knowledge transfer in the absence of the original dataset. A direct idea is to use other datasets as a substitute [27], [28], or synthesize a dataset to implement KD. Based on whether a generator is adopted, we categorize approaches using synthesized data into two classes: optimization-based and generator-based approaches.

Optimization-based approaches directly optimize a random initialized sample to make it suitable for KD. Lopes *et al.* [29] reconstructed the original dataset using metadata collected during teacher model training. Nayak *et al.* [30] produced samples by modeling the softmax space of the original dataset. Recently, statistics in batch normalization (BN) layers of the pretrained teacher model have been adopted to assist data generation [4], [31]. More realistic substitute samples can be obtained with this information. Fang *et al.* [32] further introduced contrastive learning to generate more diverse samples. However, optimization-based approaches optimize each sample separately, which means that the resource consumption is linearly related to the amount of required data. Moreover, although the scenario is assumed to be data-free, most of these approaches require the mean and variance (not those in the BN layers) of the original dataset for data generation.

Generator-based approaches adopt a generator to synthesize the substitute dataset, which is usually trained to output samples with given labels [33], [34], high teacher activation, or low prediction entropy [9]. Adversarial distillation [3], [35] is an efficient scheme that trains the generator and the student together by a min-max game. Fang *et al.* [36] and Choi *et al.* [37] combined adversarial distillation and a BN statistics constraint for data-free quantization. Liu *et al.* [10]

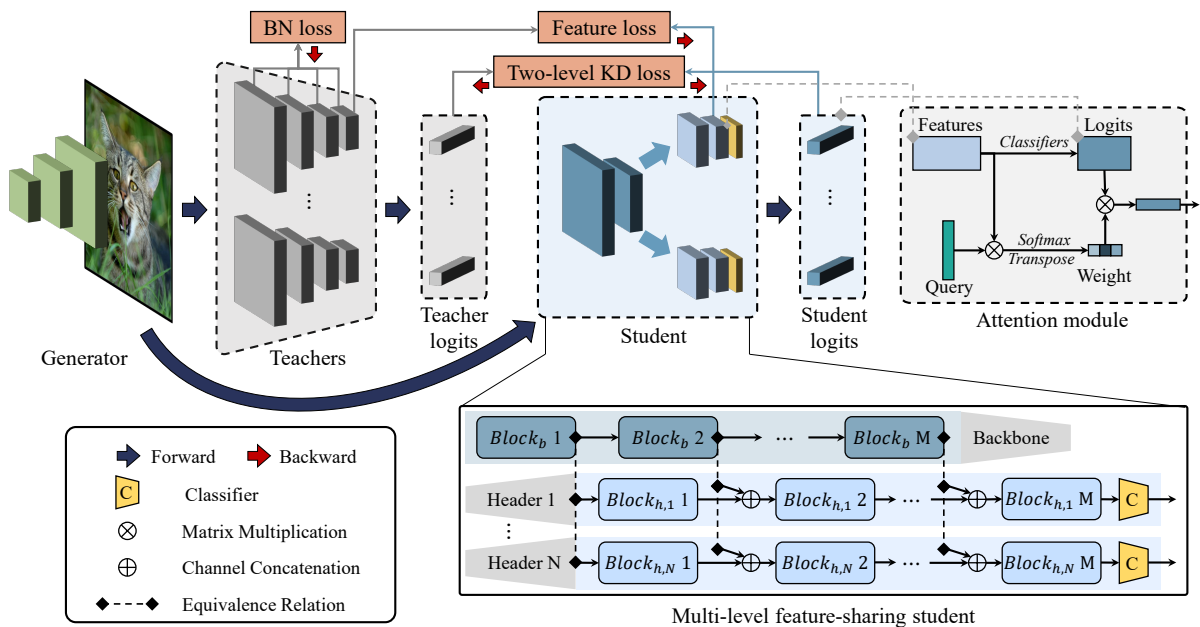


Fig. 1. Overview of the CDFKD-MFS framework. This framework consists of multiple pretrained teachers, a data generator, a multi-header student, and an attention module. In this framework, the generator produces samples to substitute the original data, and the student incorporated with multi-level feature sharing learns from teachers using these samples. When some real samples in the original dataset are available, the attention module adaptively aggregates predictions of student headers.

designed a channel relation map matching loss to further improve performance. Only a few works have studied the multi-teacher data-free KD problem. Ye *et al.* [8] designed a group-stack dual-GAN student architecture to amalgamate different knowledge. Mainly aiming at amalgamation, their student contains layers from the teacher networks, causing inefficient compression. Moreover, their student architecture is highly related to the generator, which may limit its flexibility, and the student consisting of separate trained blocks may be suboptimal. Hao *et al.* [38] used the adversarial distillation framework to achieve a multi-teacher data-free KD, but the use of a single-header student may confuse diverse knowledge in multiple teacher models. To relieve this confusion and achieve flexible knowledge transfer, we design a multi-header student and an asymmetric adversarial data-free KD algorithm, which is end-to-end.

Our previous conference work [11] performed data-free KD with multiple pretrained teacher models under a proposed framework named CDFKD. The CDFKD framework consists of multiple teacher models, a student model with multiple output headers, and a generator model. In this framework, the generator produces substitute samples, while the student learns from multiple teachers with these samples, and each teacher is assigned to a corresponding student header. Moreover, when some samples from the original dataset are available, a simple attention module is trained to aggregate predictions given by student headers.

C. Decoupling and Aggregation of Knowledge

To learn from multiple teachers and preserve the diversity, the student needs to decouple the knowledge in the teacher ensemble and learn with decoupled parameters with a branchy

network architecture. Such designs were common in early-existing [39] inference and multi-task learning. The multi-task learning problem, where early works adopted a tree-like student architecture as a solution [40], bears a closer resemblance to our work. Some later approaches decoupled the network into task-shared and task-specific parts. Misra *et al.* [41] proposed a cross-stitch unit to combine activations of the two parts. Liu *et al.* [7] proposed a multi-task attention network to achieve the decoupling of knowledge. The student designed for multi-teacher KD in our study was inspired by this work.

Aggregating predictions of the independent branches benefits from the ensemble method. We achieve aggregation with the attention mechanism, which enhances some parts of the input information while ignoring extraneous parts. The attention mechanism has been combined with DNNs to learn better feature representation [42], [43]. In addition to these simple attention mechanisms, more complex forms such as self-attention networks [44] and graph attention networks [45] have also achieved outstanding performance. In our work, we adopt key-value attention to aggregate predictions of student headers when a few real samples are given.

III. METHOD

To facilitate data-free KD to compress multiple pretrained models, we propose a framework termed CDFKD-MFS. Our framework can integrate the knowledge of multiple models into a tiny model while preserving diverse multi-model knowledge.

An overview of the proposed CDFKD-MFS framework is shown in Fig. 1. This framework contains multiple teacher models, a generator, a multi-header student model, and an

attentional aggregation module. Teacher models are pretrained on the original dataset. During data-free KD, the generator produces samples to substitute the original data. Given the generated samples, the student uses its backbone network to learn shared knowledge and its header branches to learn from specific teachers. In the student model, multi-level feature sharing is incorporated to reduce the branchy model size while providing satisfactory performance. Moreover, when some real samples from the original dataset are available, the attention module aggregates predictions of student headers to further improve performance.

In the following part of this section, we first introduce some background knowledge of data-free KD. Then we depict details of our CDFKD-MFS framework.

A. Background

Performing data-free KD requires a substitute for the original dataset. Existing approaches obtain this substitute using either generator-based or optimization-based methods. Generator-based methods train a generator network to produce substitute samples, while optimization-based methods directly optimize Gaussian noise for desired samples. Because the proposed framework is generator-based, we introduce some background knowledge and universal techniques of generator-based data-free KD, including KD, adversarial distillation, and BN statistics constraint.

KD is a technique for compressing a large pretrained model into a tiny one [2]. During training, KD uses the predictions of the pretrained model, called teacher T , to guide the learning of a tiny student model S . Denoting the parameters of the student as θ_s , the above procedure is formulated as:

$$\min_{\theta_s} \mathbf{E}_{\mathbf{x} \in \mathcal{X}} [D_{KL}(S_{\theta_s}(\mathbf{x}), T(\mathbf{x})) + \lambda D_{CE}(S_{\theta_s}(\mathbf{x}), y)],$$

where D_{KL} and D_{CE} denote the KL-divergence and cross-entropy distance, respectively, y is the ground-truth label of input sample \mathbf{x} , and λ is a hyperparameter balancing the two objectives. In practice, we found that simply setting λ to 0 has little influence on student performance, so we adopt this more straightforward KD form in our study. However, the vanilla KD approach can only work in data-available scenarios. When the original dataset is unavailable, adversarial data-free distillation is an ideal solution.

Adversarial distillation is a typical training scheme for generator-based data-free KD approaches [35]. When the original dataset is unavailable, the core problem is to find a substitute for the dataset. In adversarial distillation, this substitution is obtained by training a data generator, and the substitute data generation runs together with KD adversarially. Adversarial training is initially used in generative adversarial networks, where a generator network produces samples and a discriminator network distinguishes generated samples from real samples [46]. The two networks are trained successively in each iteration. However, the discriminator trained with the original dataset cannot be directly adopted in the case of missing data. Adversarial distillation regards the teacher and the student together as a proxy discriminator and trains it with

a negative KL divergence loss. Mathematically, the adversarial distillation procedure can be represented as:

$$\min_{\theta_s} \max_{\theta_g} \mathbf{E}_{\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})} [D_{KL}(S_{\theta_s}(G_{\theta_g}(\mathbf{z})), T(G_{\theta_g}(\mathbf{z}))),$$

where θ_g represents the parameters of generator G and \mathbf{z} is a random vector sampled from a multivariate normal distribution. This procedure can be intuitively interpreted as follows: in each iteration, the generator produces samples on which the student cannot output predictions similar to the teacher. Simultaneously, the student learns from the teacher using these samples. When the training converges, the student can learn most teacher knowledge. Hence, the quality of the generated samples significantly influences the performance of the student, and how to generate samples containing more meaningful information plays a crucial role.

BN statistics constraint has been proven helpful in sample generation for data-free KD [4], [37]. BN layers can reduce the internal covariate shifts when data forward propagate through neural networks, which makes the training of deeper networks achievable [47]. In data-free KD, these layers can also reduce the domain shift of generated samples. Specifically, during training, BN layers of the teacher model compute the mean and variance of all the data passing through them. By constraining activations of generated samples to match the precollected mean and variance on these layers, the generator can produce samples following a distribution close to the original dataset. Indexing BN layers with i and denoting the mean and variance of the i -th BN layer as μ_i and σ_i , the BN statistic constraint is given by:

$$\min_{\theta_g} \mathbf{E}_{\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})} [\|\hat{\mu}_i(G_{\theta_g}(\mathbf{z})) - \mu_i\|_2 + \|\hat{\sigma}_i(G_{\theta_g}(\mathbf{z})) - \sigma_i\|_2],$$

where $\hat{\mu}_i(G_{\theta_g}(\mathbf{z}))$ and $\hat{\sigma}_i(G_{\theta_g}(\mathbf{z}))$ are the mean and variance of activations at the i -th BN layer when the teacher takes $G_{\theta_g}(\mathbf{z})$ as input, respectively.

Although remarkable performance was achieved, almost all previous works assumed that only one pretrained teacher model exists without considering multi-teacher scenarios. We remedy this drawback and propose the CDFKD-MFS framework, which is suitable for multi-teacher data-free KD.

B. Multi-header Student Architecture

We first introduce a multi-header network architecture, which is adopted as the student in our framework. As there are multiple pretrained teachers, how to preserve their diverse knowledge into a single model requires careful consideration. A possible solution is to attach multiple output headers to a student and let each header learn from an individual teacher. Our previous work attached output headers to the student in a tree-like manner. More specifically, all headers are directly concatenated at the end of the backbone of the student. However, this simple architecture is insufficient to learn the abundant, diverse knowledge of multiple teachers and results in performance degradation, especially on some complex datasets, as shown in Section IV-C. To remedy this drawback, we propose a multi-header student with multi-level feature sharing.

TABLE I
NETWORK ARCHITECTURE OF ONE STUDENT HEADER BLOCK.

Conv($C_{in}, C_{in}, 3$), Conv($C_{in}, C_{in}, 1$), BatchNorm, ReLU
Conv($C_{in}, C_{in}, 3$), Conv($C_{in}, C_{out}, 1$), BatchNorm, ReLU

The lower right part of Fig. 1 illustrates the details of the proposed student architecture, which consists of two components: a shared backbone network and multiple header networks. There are the same number of headers and pre-trained teachers. We denote the number by N . We split the backbone network into several blocks according to where the intermediate features are shared with the headers and assume the number of feature-sharing points is M . Each header also consists of M blocks, which correspond to blocks in the backbone. The backbone can be any off-the-shelf network architecture, e.g., ResNet [1]. As the number of headers is related to the number of pretrained teachers, headers should be lightweight enough to avoid consuming too many computing resources when performing inference. Hence, we adopt the depthwise separable convolution [48] to build lightweight headers. Table I provides the architecture of one student header block, where Conv(C_1, C_2, k) denotes a convolutional layer with a $k \times k$ kernel, which takes C_1 channels feature map as input and output C_2 channels feature map.

The student takes samples produced by the generator as inputs. During forward propagation, its backbone network performs the same as an ordinary DNN, and each header network takes multi-level intermediate features of the backbone network as input. For all headers, the input of their first block is the output of the first backbone block. Then, we take the i -th header and its j -th ($j > 1$) block as the example to explain how other header blocks work. When the forward propagation of the $(j-1)$ -th backbone block and the $(j-1)$ -th block of the i -th header is finished, their output feature maps are then concatenated at the channel dimension and fed into the j -th block. Once the inference of the last block in a header is finished, the output feature is sent to the corresponding classifier to obtain the final prediction. The predictions of all headers can be aggregated by taking the average or other ensemble approaches. Similar network architecture designs can be found in multi-task learning [7], [49], which can alleviate the interference caused by learning different tasks simultaneously. In the proposed architecture, the backbone learns the knowledge shared by all headers, and each header learns from a specific teacher. This knowledge decoupling architecture can preserve the diverse knowledge of multiple teachers in a single student network and achieve better performance in comparison with a single-header student architecture, which will be shown in Section IV-C.

C. Asymmetric Adversarial Data-free KD

After the designation of the multi-header student, we introduce the proposed multi-teacher data-free KD approach. We name this approach asymmetric adversarial data-free KD because we adopt the adversarial distillation framework, but the objective functions of the generator and the student are not

entirely the same. Specifically, for the generator, its objective function has a BN statistics constraint term, and for the student, feature supervision from teachers is introduced. Next, we first introduce the three loss terms illustrated in Fig. 1, and then introduce the objection function for training the generator and the student. Finally, the optimization is introduced.

The BN loss is a constraint acting on the mean and variance of activations at the BN layers of all pretrained teachers. It requires activations of generated samples to have statistics similar to those of samples in the original dataset. In the multi-teacher scenario, we assume there are N teachers in total and index them with n . Thus, we can extend the BN statistics constraint in Section III-A to the multi-teacher setting:

$$\mathcal{L}_{bn}(\theta_g) = \frac{1}{N} \sum_{n=1}^N \mathbf{E}_{\mathbf{z} \sim \mathcal{N}(0, \mathcal{I}), i \in \mathcal{I}_n} [\|\hat{\mu}_i(G_{\theta_g}(\mathbf{z})) - \mu_i\|_2 + \|\hat{\sigma}_i(G_{\theta_g}(\mathbf{z})) - \sigma_i\|_2], \quad (1)$$

where \mathcal{I}_n is the index set of BN layers in the n -th teacher. Because multiple teachers have learned diverse feature representations from the original dataset [5], the modified BN statistics constraint takes the representations learned by all teachers into consideration. Hence, the BN loss can assist the generator in producing more informative samples than simply using a single pretrained teacher, and these samples can help achieve better knowledge transfer.

The two-level KD loss contains two loss terms designed for knowledge distillation. As mentioned above, we set the number of output headers in the student to be equal to that of pretrained teachers to assign each header a unique teacher. Thus, we can extend the adversarial distillation framework to the multi-teacher scenario:

$$\mathcal{L}_{head}(\theta_s, \theta_g) = \frac{1}{N} \sum_{n=1}^N \mathbf{E}_{\mathbf{z} \sim \mathcal{N}(0, \mathcal{I})} \|S_{\theta_s, n}(G_{\theta_g}(\mathbf{z})) - T_n(G_{\theta_g}(\mathbf{z}))\|_1, \quad (2)$$

where $S_{\theta_s, n}(\cdot)$ denotes the output coming from the n -th header of the student. Note that we replace the KL divergence measurement in the original adversarial distillation framework with the l_1 norm measurement, which has been proven to be more effective in data-free KD [50].

Taking the average predictions from multiple trained models as the final result is a simple but effective ensemble method for DNN [5]. We want to preserve the ensemble gain of pretrained teachers in our tiny student model. Thus, we design an ensemble prediction loss term to guide the ensemble of student headers to output results that resemble those of the teacher ensemble:

$$\mathcal{L}_{ens}(\theta_s) = \mathbf{E}_{\mathbf{z} \sim \mathcal{N}(0, \mathcal{I})} \left\| \frac{1}{N} \sum_{n=1}^N (S_{\theta_s}(G(\mathbf{z})) - T_n(G(\mathbf{z}))) \right\|_1. \quad (3)$$

The header loss constitutes the adversarial target, and the ensemble loss term is only used to train the student.

The feature loss introduces supervision from intermediate layers of teachers. As shown in previous works [18], [51], [52], knowledge from intermediate layers can also guide the learning of the student. We adopt this idea and design a

Algorithm 1 Asymmetric adversarial data-free KD.

Notations: learning rates η_s, η_g ; epochs E_{\max} ; iterations in one epoch I_{\max} ; student update steps in one iteration S_{\max}

Initialize($\theta_s; \theta_g$)

for epoch **in** 1 : E_{\max} ; **do**

for iteration **in** 1 : I_{\max} ; **do**

 # Knowledge Transfer

for step **in** 1 : S_{\max} ; **do**

$z \sim \mathcal{N}(0, \mathbf{I})$; samples $\leftarrow G(z)$

 compute \mathcal{L}_s with Eqn. 2, 3, 4

$$\theta_s \leftarrow \theta_s - \eta_s \frac{\partial \mathcal{L}_s}{\partial \theta_s}$$

end for

 # Discrepancy Maximization

$z \sim \mathcal{N}(0, \mathbf{I})$; samples $\leftarrow G(z)$

 compute \mathcal{L}_g with Eqn. 1, 2

$$\theta_g \leftarrow \theta_g - \eta_g \frac{\partial \mathcal{L}_g}{\partial \theta_g}$$

end for

 decay learning rates η_s, η_g

end for

feature supervision term. For simplicity and flexibility, this supervision is only imposed on the last feature extraction layer of the teachers and the student, i.e., layers whose outputs are fed into classifiers. When a model takes \mathbf{x} as input, we denote the output of these layers as $f_{t,n}(\mathbf{x})$ and $f_{s,n}(\mathbf{x})$, where t and s denote teachers and student headers, respectively, and n indexes them. Hence, we can write the feature loss term as:

$$\mathcal{L}_{feat}(\theta_s) = \frac{1}{N} \sum_{n=1}^N \mathbf{E}_{z \sim \mathcal{N}(0, \mathbf{I})} \|f_{s,n}(G(z)) - f_{t,n}(G(z))\|_1. \quad (4)$$

The objective functions of the generator and the student are composed of the above introduced loss terms. They are formulated as:

$$\begin{aligned} \mathcal{L}_s(\theta_s) &= \mathcal{L}_{head}(\theta_s, \theta_g) + \alpha \mathcal{L}_{ens}(\theta_s) + \beta \mathcal{L}_{feat}(\theta_s), \\ \mathcal{L}_g(\theta_g) &= -\mathcal{L}_{head}(\theta_s, \theta_g) + \gamma \mathcal{L}_{bn}(\theta_g), \end{aligned} \quad (5)$$

where α, β , and γ are balancing hyperparameters. The asymmetry in our adversarial data-free KD approach is apparent from the above objective functions. Compared with symmetric adversarial distillation, our approach introduces model-specific optimization objectives and can achieve better performance, which will be shown in our experiments.

The optimization of our asymmetric adversarial data-free KD resembles that in the original adversarial distillation approach, which has a two-stage form. We summarize it in Algorithm 1. We adopt the setting that trains the student more frequently than the generator to keep the adversarial process balanced [3]. Moreover, during the whole training process, we let BN layers in the generator and the student keep collecting the mean and variance of all generated samples. This setting

Algorithm 2 Attention-based prediction aggregation.

Notations: learning rates η_q

Initialize(q)

while not converged; **do**

for (x, y) **in** (\mathcal{X}, \mathcal{Y}); **do**

$z \sim \mathcal{N}(0, \mathbf{I}), \theta \sim \text{Beta}(1, 1)$

$\hat{\mathbf{x}} \leftarrow \theta \mathbf{x} + (1 - \theta)G(z)$

 compute label y' of $G(z)$ with Eqn. 9

 compute \mathcal{L}_{attn} with Eqn. 7

$$q \leftarrow q - \eta_q \frac{\partial \mathcal{L}_{attn}}{\partial q}$$

end for

end while

makes the learning of the generator and the student easy and improves the stability of the adversarial training process because they can track the activation statistics of each other.

D. Attention-based Prediction Aggregation

When some real samples are available, we use an attention module to aggregate predictions of student headers. The attention module takes a simple key-value attention form, and a Mixup [53] of real and generated samples is adopted for training the model.

The architecture of the attention-based prediction aggregation module is illustrated in the upper right corner of Fig. 1. The key vectors are composed of the features before classifiers, and the value vectors consist of logits. By computing the inner product between the key vectors and a learnable query vector and processing the result with a softmax function, we can obtain a weight vector for computing the weighted sum of all logit vectors. Then, the weighted logit vector is regarded as the final prediction. To formulate the above process, we represent key and value vectors by matrices:

$$\begin{aligned} \mathbf{K} &= [f_{s,1}^T(\mathbf{x}), f_{s,2}^T(\mathbf{x}), \dots, f_{s,N}^T(\mathbf{x})] \in \mathbb{R}^{N \times D}, \\ \mathbf{V} &= [S_{\theta_{s,1}}^T(\mathbf{x}), S_{\theta_{s,2}}^T(\mathbf{x}), \dots, S_{\theta_{s,N}}^T(\mathbf{x})] \in \mathbb{R}^{N \times C}, \end{aligned}$$

where D is the dimension of the feature vectors, C is the dimension of the logit vectors, i.e., the number of classes, and \mathbf{x} is an arbitrary input. Denoting the query vector as $\mathbf{q} \in \mathbb{R}^{D \times 1}$, we can then write the weighted result $\hat{S}_{\theta_s}(\mathbf{x})$ as:

$$\begin{aligned} \hat{S}_{\theta_s}(\mathbf{x}) &= \mathbf{V}^T \mathbf{w} \\ &= \mathbf{V}^T \text{softmax}\left(\frac{\mathbf{K} \mathbf{q}}{\sqrt{D}}\right), \end{aligned} \quad (6)$$

where \mathbf{w} is the weight vector.

Furthermore, considering that there may exist a domain gap between real and generated samples, which harms the aggregated result, we use a mix of real and generated samples

to train the query vector, resembling the Mixup approach [53]. The loss function is defined as:

$$\mathcal{L}_{attn}(\mathbf{q}) = \mathbf{E}_{(\mathbf{x}, y), \mathbf{z}, \theta} [\theta \mathcal{D}_{CE}(\hat{S}_{\theta_s}(\hat{\mathbf{x}}), y) + (1 - \theta) \mathcal{D}_{CE}(\hat{S}_{\theta_s}(\hat{\mathbf{x}}), y')], \quad (7)$$

$$\hat{\mathbf{x}} = \theta \mathbf{x} + (1 - \theta) G(\mathbf{z}), \quad (8)$$

$$y' = \arg \max \frac{1}{N} \sum_{n=1}^N T_n(G(\mathbf{z})), \quad (9)$$

$$(\mathbf{x}, y) \in (\hat{\mathcal{X}}, \hat{\mathcal{Y}}), \mathbf{z} \sim \mathcal{N}(0, \mathbf{I}), \theta \sim \text{Beta}(1, 1),$$

where $\hat{\mathbf{x}}$ is a mixed sample, y' is the label of $G(\mathbf{z})$ given by the average of teachers, $(\hat{\mathcal{X}}, \hat{\mathcal{Y}})$ is an available subset of the original dataset, and θ is a scalar sampled from a beta distribution, which controls the mixing ratio of real and generated samples. We summarize the training procedure in Algorithm 2.

IV. EXPERIMENTS

A. Datasets and Baselines

We evaluate the proposed CDFKD-MFS framework on three public computer vision datasets: CIFAR-100 [12], Caltech-101 [13], mini-ImageNet [14]. The details of each dataset are listed as follows:

- **CIFAR-100.** CIFAR-100 consists of 60000 images that belong to 100 classes uniformly. Each image is of size 32×32 . In each class, there are 500 images for training and 100 images for testing.
- **Caltech-101.** Caltech-101 consists of images belonging to 101 categories, whose sizes are approximately 300×200 . The number of images in each class ranges from 40 to 800. We randomly sample 80% of the dataset as the training set, and the remaining samples are used as the testing set. Images in this dataset are resized to 128×128 .
- **mini-ImageNet.** mini-ImageNet is designed for few-shot learning, which is a subset of the ImageNet1k dataset [54]. This dataset consists of 60000 images from 100 classes. We resplit this dataset for classification, where 80% of images sampled in each class are used for training, and the remaining 20% of images are used for testing. All images are resized to 224×224 in our experiment.

We compare the CDFKD-MFS framework with both data-required and data-free counterparts. We only select KD [2] as the data-required counterpart for comparison. Data-free approaches are either generator-based or optimization-based. In addition to the method in our conference work, we select 3 generator-based methods and 2 optimization methods as the baselines, which are listed as follows:

- **CDFKD [11].** The approach in our previous conference work for multi-teacher data-free KD, which is generator-based. We slightly modify this approach by replacing the outdated data generation module with an adversarial-trained generator.
- **DFED [38].** A generator-based approach for multi-teacher data-free KD that uses adversarial distillation and a BN statistics constraint.

TABLE II
PART OF THE TRAINING SETTINGS FOR ASYMMETRIC ADVERSARIAL DATA-FREE KD.

	CIFAR-100	Caltech-101	mini-ImageNet
E_{\max}	300	300	100
I_{\max}	50	50	750
S_{\max}	5	5	5
Batch size	256	64	64
α	5	5	1
β	0.2	0.2	0.05
γ	0.1	0.1	0.1
lr milestone	[100, 200]	[100, 200]	[30, 60, 90]

- **DFAD [3].** A generator-based approach that uses only adversarial distillation.
- **DFQ [37].** A generator-based approach that adopts adversarial distillation, a BN statistics constraint, and two extra entropy constraints on the sample category.
- **ADI [4].** An optimization-based approach that obtains substitute samples in advance and then conducts knowledge transfer with these samples.
- **CMI [32].** An optimization-based approach that adopts contrastive learning [55] for diverse sample generation.

B. Implementation Details

Models. There are three models in the proposed framework: pretrained teachers, a student, and a generator. Similar to existing data-free approaches, we adopt ResNet-34 [1] as the teacher network architecture. The number of teachers is set to 3. The only difference between these teachers is their random initialized parameters at the beginning of training, which is proven to be sufficient to obtain diverse teachers [5]. The backbone of the student is a ResNet-18 network. We set one feature-sharing point after each residual block group in the backbone, where the total number is 4. The number of headers in the student is equal to that of teachers, i.e., 3. For all header blocks, the number of their output channels equals the channel number of the shared features with which they are concatenated. The architecture of the generator comes from DCGAN [56]. For CIFAR-100, the generator adopts nearest-neighbor upsampling layers with a scale factor equal to 2 to achieve upsampling. The output block is a Tanh activation function followed by a BN layer. For Caltech-101 and mini-ImageNet, the generator adopts transposed convolutional layers with a kernel size equal to 4 for upsampling. Output block of the Caltech-101 generator is only a Tanh activation function, while the block of the mini-ImageNet generator is the same as that of the CIFAR-100 generator. For all datasets, we set the dimension of the normal noise as 256.

Training. Settings for conducting asymmetric adversarial data-free KD on the different datasets are not the same. We list the different hyperparameters in Table II, and introduce the others as follows. In optimization, the student is trained through stochastic gradient descent with a learning rate η_s of 0.1, momentum of 0.9, and weight decay of 0.0005. The generator is trained with the Adam optimizer [57], where the initial learning rate $\eta_g = 0.001$, and η_s and η_g are multiplied

TABLE III
KNOWLEDGE DISTILLATION RESULTS.

Model/Method	Year	Data-required	Params(M)	CIFAR-100		Caltech-101		mini-ImageNet	
				FLOPs(G)	Acc(%)	FLOPs(G)	Acc(%)	FLOPs(G)	Acc(%)
Teacher	-	✓	21.3	1.16	78.07	1.20	78.54	3.67	78.13
Teacher ensemble	-	✓	64.0	3.48	80.40	3.60	80.83	11.01	80.99
Student	-	✓	11.2	0.56	77.42	0.59	77.40	1.82	77.62
KD [2]	2015	✓	11.2	0.56	77.62	0.59	77.76	1.82	77.76
DFAD [3]	2019	×	11.2	0.56	68.35	0.59	70.60	1.82	56.51
DFQ [37]	2020	×	11.2	0.56	76.14	0.59	77.80	1.82	72.94
ADI [4]	2020	×	11.2	0.56	51.02	0.59	23.54	1.82	32.84
CMI [32]	2021	×	11.2	0.56	73.05	0.59	71.39	1.82	72.87
DFED [38]	2021	×	11.2	0.56	75.87	0.59	78.64	1.82	74.05
CDFKD [11]	2021	×	20.8	0.71	66.47	0.75	78.04	2.28	59.53
CDFKD-MFS(w/o attention)	-	×	18.2	0.64	77.69	0.68	79.59	2.09	76.94
CDFKD-MFS(w attention)	-	✓(10%)	18.2	0.64	77.81	0.68	79.82	2.09	77.04

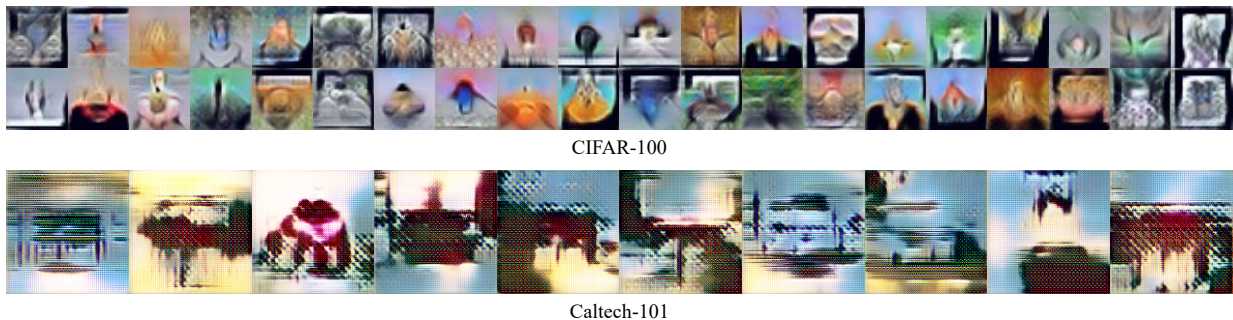


Fig. 2. Visualization of images generated on CIFAR-100 (top) and Caltech-101 (bottom). These abstract samples are suitable for knowledge transfer because each one of them contains information from more than one category.

by 0.1 when the running epoch is a member of the lr milestone array. Furthermore, we use a batch size of 128 and the AdamW optimizer [58] with a learning rate $\eta_q = 0.01$ and weight decay 0.0001 to train the query vector in the attention-based aggregation module.

All baselines are run with their official released code, except DFQ [37], which we adopt the implementation in CMI [32]. For optimization-based approaches, i.e., ADI [4] and CMI [32], we generate images whose number is equal to the original dataset in advance and then implement the knowledge transfer. We run each setting 3 times.

C. Results on CIFAR-100 and Caltech-101

We first conduct experiments on CIFAR-100 and Caltech-101, which are usually used by previous data-free counterparts. The results are summarized in Table III.

Result. We first analyze the results of CIFAR-100. Given some pretrained teacher models, the ensemble is an effective method to improve performance. By simply taking their average to obtain the new prediction, we can achieve a more than 2% accuracy improvement. Although the ensemble method performs well, it requires tremendous computing resources, which is unaffordable for edge devices. Direct training of a student model or using KD to transfer the knowledge of pretrained teachers into a student model can result in a tiny model applying to deployment. However, when the original dataset is unavailable, data-free KD approaches are required. The results show that all data-free methods can

transfer the knowledge successfully, and the proposed method performs best among them. Furthermore, when some original data are available, the performance can be further improved with the attention-based aggregation module. Our method can achieve at least 1.67% performance improvement compared with the data-free baselines. Our approach also outperforms the vanilla KD method by 0.19% because multiple pretrained teachers can provide more knowledge than a single teacher. On Caltech-101, our framework can achieve a 1.18% performance improvement over other data-free counterparts.

Discussion. According to the results of baselines with unmodified student architecture, a linear single-header student seems insufficient to preserve the diverse knowledge of multiple teachers because there is a noticeable performance gap between itself and the teachers. Hence, decoupling the knowledge by a multi-header student to bridge the gap is reasonable. Our previous CDFKD framework adopts this idea but seems ineffective. Considering that this unsatisfactory result may come from the inappropriately designed header architecture, we redesign a multi-header student with a multi-level feature-sharing structure, and it performs well. We find that the proposed architecture brings more parameters and FLOPs requirement, but different from the apparent increase in the number of parameters, owing to the depthwise separable convolution, the FLOPs requirement increases slightly, whose rate is from 10% to 20%. Compared with cheap storage, computing resources usually constrain the deployment of DNNs more. With little extra computing resources cost, our method

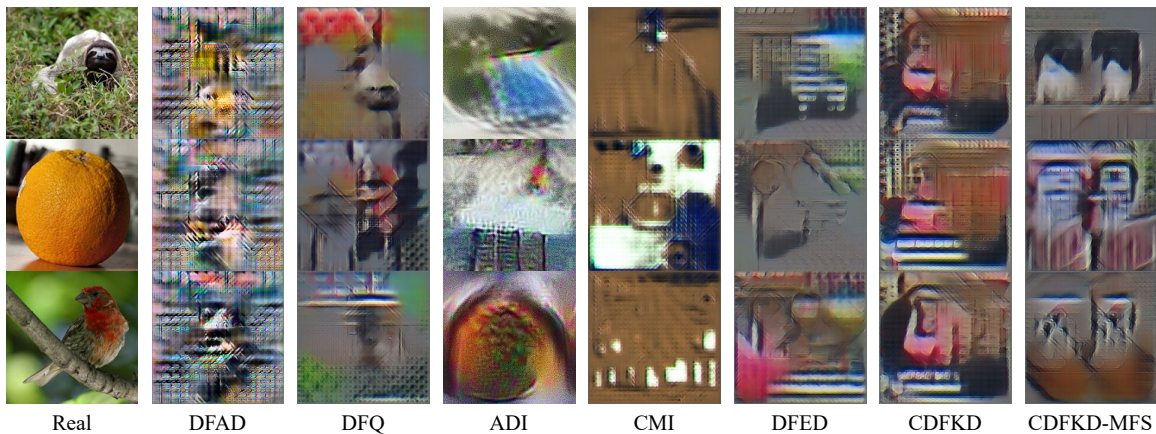


Fig. 3. Comparison of images generated on mini-ImageNet. Noting that samples are randomly selected, there is no guarantee that each row belongs to the same category. Horizontally symmetric samples generated by our method are more informative and can assist the student in learning the horizontal flip invariance of teachers.

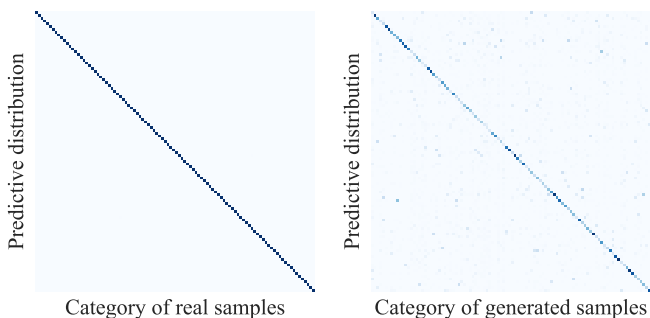


Fig. 4. Confusion matrices of real samples (left) and generated samples (right) on CIFAR-100. Generated samples contain more cross-category information than real samples.

can train a student with performance comparable to a single teacher in a data-free manner.

We also note that the proposed method even outperforms a single pretrained teacher on Caltech-101. We conjecture that this is because Caltech-101 is a simple dataset with fewer samples, so each teacher fits the training set by directly remembering most of the samples and learns fewer universal feature representations [5]. The insufficiently learned representations constrain their performance on the testing set. When multiple pretrained teachers are used, their diverse representations are aggregated into a single student model. Hence, the student learns more feature representations and performs better than a single teacher.

Sample visualization. We provide some generated samples of the proposed method in Fig. 2. These samples seem to be abstract. We conjecture that the cross-category information contained in each sample causes their confusing appearance, which means that each sample belongs to more than one category simultaneously. To verify this, we show confusion matrices of real samples and generated samples on CIFAR-100 in Fig. 4. The predictive distribution of all samples is obtained with one pretrained teacher model, and the label of each generated sample is set to the index of the max value in the

predictive distribution vector. The confusion matrices illustrate that each generated sample indeed belongs to multiple categories in the view of the teacher. According to the original KD paper [2], these samples containing cross-category information are beneficial for transferring knowledge, with which teachers can teach the student the relations among multiple categories, not only the knowledge within each individual category.

D. Results on mini-ImageNet

We then test the proposed framework on the mini-ImageNet dataset. Except for ADI, all the data-free baselines were evaluated with at most 128×128 input size in their original paper. We use the mini-ImageNet dataset with a 224×224 input size setting to test their performance on larger samples.

Result analysis. The right two columns of Table III show the results on mini-ImageNet. As we can see, the performance gap between students and teachers increases, because this dataset is complex. Compared with the data-free counterparts, the proposed method can still learn most of the knowledge and achieve at least a 2.99% performance improvement over other methods. We compare generated images from each method to demonstrate the superiority of the proposed method.

Generated sample comparison. We provide 3 samples generated by each baseline in Fig. 3. The far left column shows some samples from the original dataset. The second and third columns are samples in two generator-based single-teacher approaches: DFAD and DFQ. Without the BN statistics constraint, samples of DFAD contain more repeating texture, which is meaningless compared with DFQ. The next two columns are samples obtained with optimization-based approaches, which appear to be more like real images than the other methods, especially the samples from ADI. Nevertheless, these methods require the precollected mean and variance of the original training samples, which may not be available in some more strict data-free scenarios. Moreover, optimization-based approaches perform worse than the generator-based counterparts in our experiments. The last three columns show samples generated by the multi-teacher methods. Samples in DFED are similar to those in DFQ because DFED regards the

TABLE IV
ABLATION STUDY ON AUXILIARY LOSS TERMS.

Ensemble loss	Feature loss	Accuracy(%)
×	×	73.62
✓	×	76.84
×	✓	77.24
✓	✓	77.69

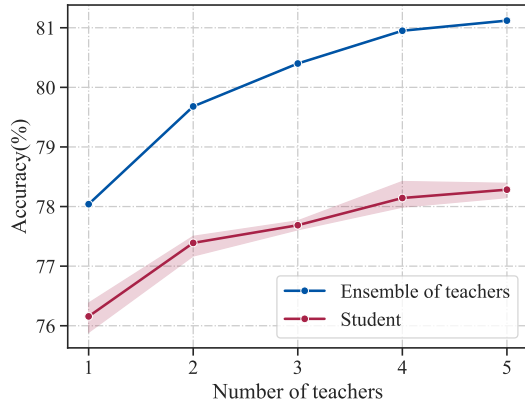


Fig. 5. Impact of teacher number. As the teacher number increases, the performance of the student increases as well.

average of multiple teachers as one teacher. However, samples in CDFKD have a similar appearance, which implies that mode collapse [46] occurred in the generator. Mode collapse means that a generator can only output a single image with a few variations, which is a common problem in GANs. We conjecture that the tree-like student confused the knowledge in the backbone during backpropagation, which hinders more useful information flowing into the generator. Samples generated in the proposed CDFKD-MFS framework appear different from the other methods. More specifically, they seem horizontally symmetric because random horizontal flipping is adopted as one of the data augmentation schemes when training teachers. This symmetry means that the generator can produce samples containing more teacher knowledge, and hence, the student achieves better performance.

E. Ablation Study

To understand each part of the proposed framework, we design ablation studies, where the framework is evaluated on CIFAR-100 with no available sample. We study the effectiveness of two auxiliary loss terms and the impact of teacher number.

Auxiliary loss. There are three auxiliary loss terms for training the student and the teacher. They are asymmetric parts in the objective functions (5), controlled by hyperparameters α , β , and γ . Because the BN statistics constraint has been proven effective, we only consider the ensemble loss (3) and the feature loss (4) here. Table IV shows the evaluation result. Each auxiliary loss can help improve the performance, and the best result is achieved when both terms are adopted. The improvement is 4.07% when both losses are used.

Impact of teacher number. The ensemble method is an effective approach to improve model performance. We vary the number of teachers in the ensemble from 1 to 5 and evaluate how the student performs. The result is shown in Fig. 5. Not surprisingly, the accuracy of the student increased with the number of teachers. Moreover, as the teacher number increases, the redundant part of the knowledge increases. Hence, the performance growth rate decreases.

F. Results on More Lightweight Models

DFED has shown that it is challenging to perform data-free KD on small models, e.g., wide residual networks (WRNs) [59], because it results in remarkable accuracy loss. To further verify the effectiveness, we compare the proposed framework with previous works on CIFAR-100 with WRN models. All hyperparameters are the same as those in Section IV-C, and the attention module is not used.

The results are reported in Table V. We first take a ResNet-18 network as the teacher and a WRN-40-2 network as the student, which contains much fewer parameters than the teacher. Compared with data-required methods, there is an apparent performance degradation when applying data-free methods, larger than 5%. We then exchange the teacher and the student and observe a smaller performance gap of only approximately 2%. Hence, we conjecture that performing data-free KD gets more challenging as the student becomes smaller. We take two WRN models with fewer parameters as the student and a WRN-40-2 model as the teacher. The results follow our conjecture, where the performance gap between data-free and data-required methods becomes more significant when the student becomes smaller. However, our framework can significantly narrow the performance gap. When we take the WRN-16-2 model as the student, the proposed method achieves an 8.32% performance improvement compared with the data-free counterparts, and it achieves a 16.51% performance improvement when we use a smaller WRN-16-1 network as the student model.

G. Attention Map Comparison

We further compare attention maps generated by different teachers and student headers on the same sample to study what the student has learned. Attention maps are obtained on the last residual block of each model with GradCAM [60]. Samples for comparison come from the mini-ImageNet dataset.

Attention maps are shown in Fig. 6. We select four different samples according to the number of objects in the sample and the similarity of teacher attention maps. As we can see, when there is only one object in the sample, the student learns from teachers well. Each student header gives attention maps similar to its corresponding teacher (Fig. 6a, 6b). However, when multiple objects exist in a single sample, the result is different. The student either merges distinct teacher attention regions (Fig. 6c), or confuses attention from different teachers (Fig. 6d). Hence, we conjecture that samples containing multiple objects belonging to the label category are the main obstacle for data-free KD. In our future work, we will study how to guide a student to learn from teachers on such samples.

TABLE V
KNOWLEDGE DISTILLATION RESULTS ON THE CIFAR-100 DATASET WITH WRN MODELS.

Teacher			Student			Accuracy(%)					
Model	Params(M)	FLOPs(M)	Model	Params(M)	FLOPs(M)	T.*	S.*	KD*	DFAD	DFED	CDFKD-MFS
ResNet-18	11.22	557	WRN-40-2	2.26/2.71	329/344	77.42	75.77	75.90	45.05	69.15	70.26
WRN-40-2	2.26	329	ResNet-18	11.22/18.16	557/639	75.77	77.42	77.91	49.86	74.58	75.80
WRN-40-2	2.26	329	WRN-16-2	0.70/1.16	102/116	75.77	71.67	72.02	30.11	60.76	69.08
WRN-40-2	2.26	329	WRN-16-1	0.18/0.31	27/31	75.77	65.86	65.98	10.54	41.06	57.57

The left and the right of a slash (/) are data of a linear student and a multi-header student, respectively. Data-required methods are marked by an asterisk (*).

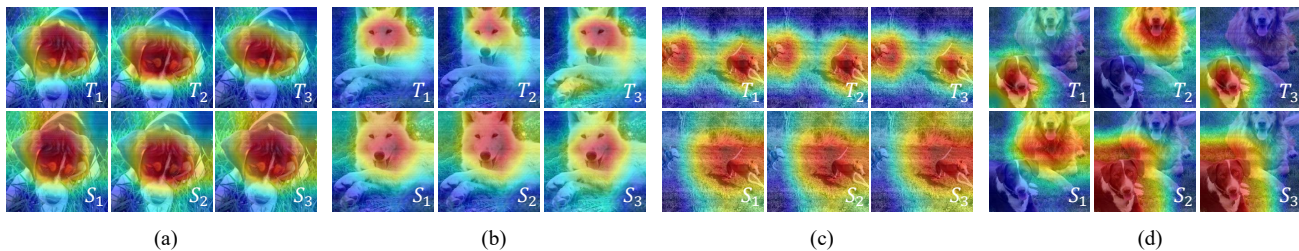


Fig. 6. Teacher-student attention map comparison: (a) single object with similar teacher attention maps; (b) single object with different teacher attention maps; (c) multiple objects with similar teacher attention maps; and (d) multiple objects with different teacher attention maps. The lower right corner in each image shows which teacher or student header generates the attention map.

V. CONCLUSION

In this paper, we propose a framework termed CDFKD-MFS for multi-teacher data-free KD. The framework consists of a multi-level feature-sharing multi-header student model architecture, an asymmetric adversarial data-free KD algorithm, and an attention-based prediction aggregation module. We conduct extensive experiments to evaluate our framework. The results show that in addition to the widely used CIFAR-100 dataset and Caltech-101 dataset, the proposed framework also outperforms other data-free counterparts on the mini-ImageNet dataset under the more challenging 224×224 input size. When some real samples are available, our method achieves better performance by using the attention mechanism to adaptively aggregate predictions of the multi-header student. Moreover, when the student model is extremely small, the proposed framework can still achieve satisfactory performance, which is difficult for existing works. In the future, we will explore how to teach a student in data-free KD to perform like the teachers on samples containing multiple objects and further improve performance.

REFERENCES

- [1] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, Las Vegas, NV, USA, June 27-30, 2016, pp. 770–778.
- [2] G. E. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *CoRR*, vol. abs/1503.02531, 2015.
- [3] G. Fang, J. Song, C. Shen, X. Wang, D. Chen, and M. Song, "Data-free adversarial distillation," *CoRR*, vol. abs/1912.11006, 2019.
- [4] H. Yin, P. Molchanov, J. M. Alvarez, Z. Li, A. Mallya, D. Hoiem, N. K. Jha, and J. Kautz, "Dreaming to distill: Data-free knowledge transfer via deepinversion," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, Seattle, WA, USA, 2020, pp. 8712–8721.
- [5] Z. Allen-Zhu and Y. Li, "Towards understanding ensemble, knowledge distillation and self-distillation in deep learning," *CoRR*, vol. abs/2012.09816, 2020.
- [6] L. Tran, B. S. Veeling, K. Roth, J. Swiatkowski, J. V. Dillon, J. Snoek, S. Mandt, T. Salimans, S. Nowozin, and R. Jenatton, "Hydra: Preserving ensemble diversity for model distillation," *CoRR*, vol. abs/2001.04694, 2020.
- [7] S. Liu, E. Johns, and A. J. Davison, "End-to-end multi-task learning with attention," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, Long Beach, CA, USA, June 16-20, 2019, pp. 1871–1880.
- [8] J. Ye, Y. Ji, X. Wang, X. Gao, and M. Song, "Data-free knowledge amalgamation via group-stack dual-gan," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, Seattle, WA, USA, June 13-19, 2020, pp. 12513–12522.
- [9] H. Chen, Y. Wang, C. Xu, Z. Yang, C. Liu, B. Shi, C. Xu, C. Xu, and Q. Tian, "Data-free learning of student networks," in *IEEE/CVF International Conference on Computer Vision, ICCV*, Seoul, Korea (South), October 27 - November 2, 2019, pp. 3513–3521.
- [10] Y. Liu, W. Zhang, and J. Wang, "Zero-shot adversarial quantization," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, virtual, June 19-25, 2021, pp. 1512–1521.
- [11] Z. Hao, Y. Luo, Z. Wang, H. Hu, and J. An, "Model compression via collaborative data-free knowledge distillation for edge intelligence," in *IEEE International Conference on Multimedia and Expo, ICME*. IEEE, 2021, pp. 1–6.
- [12] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [13] F. Li, R. Fergus, and P. Perona, "One-shot learning of object categories," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 28, no. 4, pp. 594–611, 2006.
- [14] O. Vinyals, C. Blundell, T. Lillicrap, K. Kavukcuoglu, and D. Wierstra, "Matching networks for one shot learning," in *Advances in Neural Information Processing Systems, NeurIPS*, Barcelona, Spain, December 5-10, 2016, pp. 3630–3638.
- [15] Q. Ding, S. Wu, H. Sun, J. Guo, and S. Xia, "Adaptive regularization of labels," *CoRR*, vol. abs/1908.05474, 2019.
- [16] R. Müller, S. Kornblith, and G. E. Hinton, "When does label smoothing help?" in *Advances in Neural Information Processing Systems, NeurIPS*, Vancouver, BC, Canada, December 8-14, 2019, pp. 4696–4705.
- [17] S. Kim and H. Kim, "Transferring knowledge to smaller network with class-distance loss," in *International Conference on Learning Representations, ICLR*, Toulon, France, April 24-26, 2017.
- [18] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, "Fitnets: Hints for thin deep nets," in *International Conference on Learning Representations, ICLR*, San Diego, CA, USA, May 7-9, 2015.
- [19] X. Wang, T. Fu, S. Liao, S. Wang, Z. Lei, and T. Mei, "Exclusivity-consistency regularized knowledge distillation for face recognition," in

- European Conference on Computer Vision, ECCV*, vol. 12369, Glasgow, UK, August 23-28, 2020, pp. 325–342.
- [20] K. Xu, L. Rui, Y. Li, and L. Gu, “Feature normalized knowledge distillation for image classification,” in *European Conference on Computer Vision, ECCV*, vol. 12370, Glasgow, UK, August 23-28, 2020, pp. 664–680.
- [21] S. Zagoruyko and N. Komodakis, “Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer,” in *International Conference on Learning Representations, ICLR*, Toulon, France, April 24-26, 2017.
- [22] N. Passalis and A. Tefas, “Learning deep representations with probabilistic knowledge transfer,” in *European Conference on Computer Vision, ECCV*, vol. 11215, Munich, Germany, September 8-14, 2018, pp. 283–299.
- [23] D. Chen, J. Mei, Y. Zhang, C. Wang, Z. Wang, Y. Feng, and C. Chen, “Cross-layer distillation with semantic calibration,” in *AAAI Conference on Artificial Intelligence, AAAI*, Virtual Event, February 2-9, 2021, pp. 7028–7036.
- [24] S. You, C. Xu, C. Xu, and D. Tao, “Learning from multiple teacher networks,” in *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, Halifax, NS, Canada, August 13-17, 2017, pp. 1285–1294.
- [25] S. Park and N. Kwak, “Feature-level ensemble knowledge distillation for aggregating knowledge from multiple networks,” in *European Conference on Artificial Intelligence, ECAI*, vol. 325, Santiago de Compostela, Spain, August 29 - September 8, 2020, pp. 1411–1418.
- [26] U. Asif, J. Tang, and S. Herrer, “Ensemble knowledge distillation for learning improved and efficient networks,” in *European Conference on Artificial Intelligence, ECAI*, vol. 325, Santiago de Compostela, Spain, August 29 - September 8, 2020, pp. 953–960.
- [27] G. K. Nayak, K. R. Mopuri, and A. Chakraborty, “Effectiveness of arbitrary transfer sets for data-free knowledge distillation,” in *IEEE Winter Conference on Applications of Computer Vision, WACV*, Waikoloa, HI, USA, January 3-8, 2021, pp. 1429–1437.
- [28] H. Chen, T. Guo, C. Xu, W. Li, C. Xu, C. Xu, and Y. Wang, “Learning student networks in the wild,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, virtual, June 19-25, 2021, pp. 6428–6437.
- [29] R. G. Lopes, S. Fenu, and T. Starner, “Data-free knowledge distillation for deep neural networks,” *CoRR*, vol. abs/1710.07535, 2017.
- [30] G. K. Nayak, K. R. Mopuri, V. Shaj, V. B. Radhakrishnan, and A. Chakraborty, “Zero-shot knowledge distillation in deep networks,” in *International Conference on Machine Learning, ICML*, vol. 97, Long Beach, California, June 9-15, 2019, pp. 4743–4751.
- [31] M. Haroush, I. Hubara, E. Hoffer, and D. Soudry, “The knowledge within: Methods for data-free model compression,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, Seattle, WA, USA, June 13-19, 2020, pp. 8491–8499.
- [32] G. Fang, J. Song, X. Wang, C. Shen, X. Wang, and M. Song, “Contrastive model inversion for data-free knowledge distillation,” *CoRR*, vol. abs/2105.08584, 2021.
- [33] J. Yoo, M. Cho, T. Kim, and U. Kang, “Knowledge extraction with no observable data,” in *Advances in Neural Information Processing Systems, NeurIPS*, Vancouver, BC, Canada, December 8-14, 2019, pp. 2701–2710.
- [34] L. Luo, M. Sandler, Z. Lin, A. Zhmoginov, and A. Howard, “Large-scale generative data-free distillation,” *CoRR*, vol. abs/2012.05578, 2020.
- [35] P. Micaelli and A. J. Storkey, “Zero-shot knowledge transfer via adversarial belief matching,” in *Advances in Neural Information Processing Systems, NeurIPS*, Vancouver, BC, Canada, December 8-14, 2019, pp. 9547–9557.
- [36] S. Xu, H. Li, B. Zhuang, J. Liu, J. Cao, C. Liang, and M. Tan, “Generative low-bitwidth data free quantization,” in *European Conference on Computer Vision, ECCV*, vol. 12357, Glasgow, UK, August 23-28, 2020, pp. 1–17.
- [37] Y. Choi, J. P. Choi, M. El-Khamy, and J. Lee, “Data-free network quantization with adversarial knowledge distillation,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops*, Seattle, WA, USA, June 14-19, 2020, pp. 3047–3057.
- [38] Z. Hao, Y. Luo, H. Hu, J. An, and Y. Wen, “Data-free ensemble knowledge distillation for privacy-conscious multimedia model compression,” in *ACM Multimedia, MM*, Virtual Event, China, October 20-24, 2021, pp. 1803–1811.
- [39] S. Teerapittayanon, B. McDanel, and H. T. Kung, “Branchynet: Fast inference via early exiting from deep neural networks,” *CoRR*, vol. abs/1709.01686, 2017.
- [40] L. Han, Y. Zhang, G. Song, and K. Xie, “Encoding tree sparsity in multi-task learning: A probabilistic framework,” in *AAAI Conference on Artificial Intelligence, AAAI*, Québec City, Québec, Canada, July 27-31, 2014, pp. 1854–1860.
- [41] I. Misra, A. Shrivastava, A. Gupta, and M. Hebert, “Cross-stitch networks for multi-task learning,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, Las Vegas, NV, USA, June 27-30, 2016, pp. 3994–4003.
- [42] S. Woo, J. Park, J. Lee, and I. S. Kweon, “CBAM: convolutional block attention module,” in *European Conference on Computer Vision, ECCV*, vol. 11211, Munich, Germany, September 8-14, 2018, pp. 3–19.
- [43] J. Hu, L. Shen, and G. Sun, “Squeeze-and-excitation networks,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, Salt Lake City, UT, USA, June 18-22, 2018, pp. 7132–7141.
- [44] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems, NeurIPS*, Long Beach, CA, USA, December 4-9, 2017, pp. 5998–6008.
- [45] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, “Graph attention networks,” in *International Conference on Learning Representations, ICLR*, Vancouver, BC, Canada, April 30 - May 3, 2018.
- [46] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. C. Courville, and Y. Bengio, “Generative adversarial networks,” *CoRR*, vol. abs/1406.2661, 2014.
- [47] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” in *International Conference on Machine Learning, ICML*, vol. 37, Lille, France, 6-11 July 2015, pp. 448–456.
- [48] F. Chollet, “Xception: Deep learning with depthwise separable convolutions,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, Honolulu, HI, USA, July 21-26, 2017, pp. 1800–1807.
- [49] A. A. Rusu, N. C. Rabinowitz, G. Desjardins, H. Soyer, J. Kirkpatrick, K. Kavukcuoglu, R. Pascanu, and R. Hadsell, “Progressive neural networks,” *CoRR*, vol. abs/1606.04671, 2016.
- [50] J. Truong, P. Maini, R. J. Walls, and N. Papernot, “Data-free model extraction,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, virtual, June 19-25, 2021, pp. 4771–4780.
- [51] J. Yim, D. Joo, J. Bae, and J. Kim, “A gift from knowledge distillation: Fast optimization, network minimization and transfer learning,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, Honolulu, HI, USA, July 21-26, 2017, pp. 7130–7138.
- [52] Z. Huang and N. Wang, “Like what you like: Knowledge distill via neuron selectivity transfer,” *CoRR*, vol. abs/1707.01219, 2017.
- [53] H. Zhang, M. Cissé, Y. N. Dauphin, and D. Lopez-Paz, “mixup: Beyond empirical risk minimization,” in *International Conference on Learning Representations, ICLR*, Vancouver, BC, Canada, April 30 - May 3, 2018.
- [54] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. S. Bernstein, A. C. Berg, and L. Fei-Fei, “Imagenet large scale visual recognition challenge,” *Int. J. Comput. Vis.*, vol. 115, no. 3, pp. 211–252, 2015.
- [55] T. Chen, S. Kornblith, M. Norouzi, and G. E. Hinton, “A simple framework for contrastive learning of visual representations,” in *International Conference on Machine Learning, ICML*, vol. 119, Virtual Event, July 13-18 2020, pp. 1597–1607.
- [56] A. Radford, L. Metz, and S. Chintala, “Unsupervised representation learning with deep convolutional generative adversarial networks,” in *International Conference on Learning Representations, ICLR*, San Juan, Puerto Rico, May 2-4, 2016.
- [57] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *International Conference on Learning Representations, ICLR*, San Diego, CA, USA, May 7-9, 2015.
- [58] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” in *International Conference on Learning Representations, ICLR*, New Orleans, LA, USA, May 6-9, 2019.
- [59] S. Zagoruyko and N. Komodakis, “Wide residual networks,” in *British Machine Vision Conference, BMVC*, York, UK, September 19-22, 2016.
- [60] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, “Grad-cam: Visual explanations from deep networks via gradient-based localization,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, Venice, Italy, October 22-29, 2017, pp. 618–626.