Teacher as a Lenient Expert: Teacher-Agnostic Data-Free Knowledge Distillation

Hyunjune Shin, Dong-Wan Choi*

Department of Computer Science and Engineering, Inha University, South Korea heounjunee@gmail.com, dchoi@inha.ac.kr

Abstract

Data-free knowledge distillation (DFKD) aims to distill pretrained knowledge to a student model with the help of a generator without using original data. In such data-free scenarios, achieving stable performance of DFKD is essential due to the unavailability of validation data. Unfortunately, this paper has discovered that existing DFKD methods are quite sensitive to different teacher models, occasionally showing catastrophic failures of distillation, even when using well-trained teacher models. Our observation is that the generator in DFKD is not always guaranteed to produce precise yet diverse samples using the existing representative strategy of minimizing both class-prior and adversarial losses. Through our empirical study, we focus on the fact that class-prior not only decreases the diversity of generated samples, but also cannot completely address the problem of generating unexpectedly low-quality samples depending on teacher models. In this paper, we propose the teacher-agnostic data-free knowledge distillation (TA-DFKD) method, with the goal of more robust and stable performance regardless of teacher models. Our basic idea is to assign the teacher model a *lenient* expert role for evaluating samples, rather than a strict supervisor that enforces its class-prior on the generator. Specifically, we design a sample selection approach that takes only clean samples verified by the teacher model without imposing restrictions on the power of generating diverse samples. Through extensive experiments, we show that our method successfully achieves both robustness and training stability across various teacher models, while outperforming the existing DFKD methods.

Introduction

Knowledge distillation (KD) (Hinton, Vinyals, and Dean 2015) is a powerful compression technique that transfers the knowledge of a pretrained teacher model to a smaller student model. Typically, KD methods require data samples that are used to train the teacher model, in order to properly guide the training of the student model. However, in real-world scenarios, it is neither always possible nor desirable to assume the availability of training data. To address such practical issues, data-free knowledge distillation (DFKD) has been actively

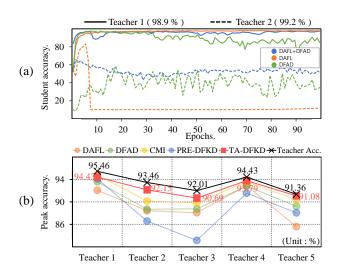


Figure 1: (a) Training curves of student models during distillation of two different well-trained teacher models on MNIST when using one or both of class-prior (DAFL) and adversarial learning (DFAD) losses. (b) Peak accuracies of student models distilled from five different teacher models on CIFAR10 when using different DFKD methods including our TA-DFKD method

studied (Yoo et al. 2019; Nayak et al. 2019), aiming to distill pretrained knowledge through the assistance of a generator, without the use of original data samples. The generator is also trained based on the teacher model to generate synthetic samples, which are intended to be replacements of the original samples in the distillation process.

A key challenge in DFKD arises from the unavailability of validation data, making it impossible to accurately evaluate the effectiveness of distillation. Therefore, it is crucial for DFKD methods to ensure the stable and robust performance no matter which teacher models are distilled. To this end, the state-of-the-art (SOTA) DFKD methods incorporate the following three components into their training loss function for the generator: *class-prior*, *adversarial*, and *representation* losses. Class-prior, initially introduced by DAFL (Chen et al. 2019), aims to generate accurate samples that can be classified by the teacher model into a specific class. On the other hand, the adversarial loss, first proposed by DFAD (Fang

^{*}Corresponding Author Copyright © 2024, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

et al. 2019) and ZSKT (Micaelli and Storkey 2019), intends to generate hard samples that maximize the output discrepancy between the teacher and student models, thereby enhancing the diversity of generated samples. Lastly, the representation loss focuses on learning feature-level information of real data with respect to the teacher model.

Unfortunately, this paper has discovered that existing DFKD methods are quite sensitive to different teacher models, occasionally showing catastrophic failures of distillation, even when using well-trained teacher models with high performance. As shown in Figure 1(a), we first focus on how two representative techniques, such as enforcing class-prior (e.g., DAFL) and minimizing adversarial loss (e.g., DFAD), can sometimes fail in distillation from two different teacher models on MNIST with the same level of accuracy achieved by the same training method. Although *Teacher 2* is slightly better in accuracy than Teacher 1, both DAFL and DFAD completely fail to distill the knowledge of Teacher 2, resulting in a large performance gap from their successful counterparts using Teacher 1. Although more and more recent proposals combine both class-prior and adversarial losses, such a mixed approach may not be a successful solution either, as shown by the failure of DAFL + DFAD in distilling *Teacher 2* in Figure 1(a). Note that there is no difference in training strategies between Teacher 1 and Teacher 2.

In our findings, this teacher-sensitive failure in DFKD occurs mainly due to a misguided generator that does not always produce precise yet diverse samples when employing the above two strategies, namely minimizing class-prior and adversarial losses. First of all, class-prior such as in DAFL is intended to improve the sample quality, but it also tends to guide the generator to focus on only easy samples. As a result, the student model can learn only a small fraction of the teacher's knowledge. In extreme cases, the resulting student model can misclassify every sample into a particular class, which is why the model distilled from Teacher 2 by DAFL keeps 0.1 accuracy (i.e., out of 10 digits) in Figure 1(a). On the other hand, the adversarial loss for the generator used in DAFD is effective to generate harder samples, which are possibly more diverse as well, but can lead to unrealistic samples that are not relevant to any of the classes of the teacher model. In order to achieve both high quality and diversity of synthetic samples, the recent works like (Fang et al. 2021; Yin et al. 2020; Binici et al. 2022a,b; Li et al. 2023) combine both techniques. However, depending on teacher models, we find that they are not guaranteed to find a sweet spot between two conflicting losses, one for precision and the other for diversity, and consequently suffer from the generation of unexpectedly low-quality samples. As demonstrated in Figure 1(b), none of the SOTA DFKD methods show a satisfactory level of robustness across 5 different pretrained models on CIFAR-10.

In this paper, we revisit the necessity of class-prior, which has been believed crucial by most SOTA methods, and focus on its drawback, namely enforcing the teacher's strict restriction to the generator. In our analysis, we find that a generator can freely generate more diverse samples when it is trained without class-prior. Moreover, despite the attempts of class-prior to enhance the sample quality, our observation

	DAFL	DFAD	ADI	CMI	PRE-DFKD	TA-DFKD
Cls.	√		✓	✓	✓	
Adv.		✓	✓	✓	✓	✓
Rep.	√ activ.		√ BNS	√ BNS	√ activ.	√ BNS

Table 1: Summary of the existing DFKD methods, DAFL (Chen et al. 2019), DFAD (Fang et al. 2019), ADI (Yin et al. 2020), CMI (Fang et al. 2021), PRE-DFKD (Binici et al. 2022a) and TA-DFKD (ours), in terms of using three major components, class-prior, adversarial, and representation losses.

reveals that relying only on the class-prior loss still allows the generator to produce low-quality samples, even without the adversarial loss.

Based on these observations, we propose the *teacher-agnostic data-free knowledge distillation* (TA-DFKD) method that assigns the teacher model a *lenient* expert role, namely removing the class-prior restriction for the generator to explore larger area in the sample space for achieving higher diversity. At the same time, in pursuit of high precision of synthetic samples, TA-DFKD utilizes the teacher model as an expert who can evaluate the quality of synthetic samples, thereby discards unexpectedly low-quality samples. Inspired by the existing works (Song et al. 2020) on learning from noisy labels, we design a sample selection method that takes only generated samples whose labels are confirmed to be sufficiently precise by the teacher model, using the Gaussian Mixture Model.

As observed in Figure 1(b), our TA-DFKD method demonstrates a highest level of teacher-agnostic robustness by consistently achieving the best accuracy close to those of teacher models. This trend is also observed in our extensive experimental results, where TA-DFKD manages to achieve both the robustness across various teacher models and stability at converging time of the distillation process, outperforming the existing DFKD methods.

Related Works

Data-Free Knowledge Distillation In data-free knowledge distillation (DFKD), given only a pretrained teacher model without any real or meta data, our focus is on how to generate synthetic samples that can be used to effectively transfer the teacher's knowledge to a target student model. There are two initial strategies to this end, optimizing random noisy images themselves (Nayak et al. 2019) or employing a generator extracted from a pretrained model (Yoo et al. 2019; Chen et al. 2019). Since the former is more computationally expensive (Nayak et al. 2019; Yin et al. 2020), recent studies have primarily focused on the latter approach, where the main issue is to train the generator only using the teacher model. Except for KegNet (Yoo et al. 2019), most DFKD methods (Chen et al. 2019; Fang et al. 2019; Micaelli and Storkey 2019; Binici et al. 2022b,a; Do et al. 2022; Li et al. 2023) adopt a one-phase distillation scheme such that the generator and the student are simultaneously trained from scratch, while freezing the teacher network. This enables a progressive transfer of the teacher's knowledge using the generator being trained. To generate more effective samples, three types of loss terms are mainly leveraged for training the generator: *class-prior*, *adversarial*, and *representation* losses, as described in the previous section. Table 1 provides a summary of which losses are employed in existing DFKD methods.

DAFL (Chen et al. 2019) first exploits class-prior that enforces the generator to produce samples that are precise enough to be well predicted by the teacher. It also proposes a representation loss, referred to as *activation*, which aims to maximize activation values of the feature maps. Another representation loss, introduced in ADI (Yin et al. 2020) and called BNS, constrains the statistics of batch normalization layers stored in the teacher model. DFAD (Fang et al. 2019) and ZSKT (Micaelli and Storkey 2019) adopt an adversarial learning strategy inspired by GAN (Goodfellow et al. 2014), aiming to generate more challenging samples that maximize disagreement between the teacher and student models. This approach encourages the student model to learn diverse knowledge from the teacher model, but may lead to unrealistic samples that belongs to none of teacher's categories.

Consequently, recent studies have attempted to combine class-prior, adversarial and representation losses, aiming of generating precise and diverse samples, while proposing their additional techniques to further enhance performance. CMI (Fang et al. 2021) suggests using contrastive learning to increase the diversity of generated samples. CuDFKD (Li et al. 2023) and AdaDFQ (Qian et al. 2023) propose adaptive learning so that the student model can progressively learn the teacher's knowledge, and ABD (Hong et al. 2023) deals with a scenario with untrustworthy teacher models. Furthermore, MB/PRE-DFKD (Binici et al. 2022b,a) pay attention to undesirable forgetting in the student model caused by adversarial learning, as seen in the training curve of DFAD with *Teacher 1* in Figure 1(a). To prevent this forgetting phenomenon, MB/PRE-DFKD (Binici et al. 2022b,a) propose the use of a memory bank or an extra generative model. With the same goal, MAD (Do et al. 2022) suggests employing exponential moving average for generator updates, while META-DFKD (Patel, Mopuri, and Qiu 2023) incorporates meta-learning into the generator training process. Despite some synergy effects observed in these DFKD methods that combine the three losses, none of them achieve a satisfactory level of robustness and stability across different teacher models, as revealed in our experimental results.

Learning from Noisy Labels Unlike popular benchmark datasets assuming always correct labels in deep neural networks (DNNs), data labeling in practice can be highly prone to errors, leading to noisy labels. To address this issue, there has been a branch of works, called learning from noisy labels (LNL) (Song et al. 2020), which focuses on preventing a DNN from overfitting to data with noisy labels. A representative approach is sample selection that identifies clean samples by modeling the difference between clean ones and those with noisy labels. A simple policy can be taking sam-

ples with smaller loss values. More advanced strategies include using a pretrained model (Jiang et al. 2018) and training dual models (Malach and Shalev-Shwartz 2017; Han et al. 2018; Yu et al. 2019; Li, Socher, and Hoi 2020) to make a better decision on clean samples. Our method is inspired by these sample selection strategies in LNL even though the DFKD problem itself is not directly related to identifying noisy labels. Specifically, we employ the Gaussian Mixture Model introduced by DivideMix (Li, Socher, and Hoi 2020), as it aligns well with our objective of selecting high-quality samples with respect to the pretrained teacher model.

Methodology

Framework of Generator-Based DFKD

In the standard generator-based DFKD framework, we consider the following three networks: a pretrained teacher model θ_T , a student model θ_S , and a generator θ_G . The ultimate goal of DFKD is the same as in the normal KD, that is, transferring the knowledge of the teacher model to the student model. Instead of real data, however, DFKD uses the generator to generate a fake sample $\hat{x} = \theta_G(z)$ with some random vector $z \sim p_z(z)$, and feeds these synthetic samples to θ_T and θ_S for minimizing the following distillation loss:

$$\mathcal{L}_{KD} = \mathbb{E}_{z \sim p_z(z)} [\mathcal{D}(\theta_T(\theta_G(z)), \ \theta_S(\theta_G(z)))], \quad (1)$$

where $\mathcal{D}(\cdot,\cdot)$ is the distance between the outputs of two models and $p_z(z)$ is usually $\mathcal{N}(0,1)$. The most challenging issue here is how to define an effective loss function \mathcal{L}_G to train θ_G without using any real data. To this end, most DFKD methods employ a mixed loss function as follows:

$$\mathcal{L}_G = \alpha \mathcal{L}_{Cls} + \beta \mathcal{L}_{Adv} + \gamma \mathcal{L}_{Rep}, \tag{2}$$

where \mathcal{L}_{Cls} is the class-prior loss, \mathcal{L}_{Adv} is the adversarial loss, and \mathcal{L}_{Rep} is the representation loss. Given \mathcal{L}_{KD} and \mathcal{L}_{G} , while freezing θ_{T} , the final goal of DFKD is to simultaneously train θ_{S} and θ_{G} with the following objective functions: $\min_{\theta_{S}} \mathcal{L}_{KD}$ and $\min_{\theta_{G}} \mathcal{L}_{G}$.

Our Findings. In this work, we argue that the generator is not always guaranteed to synthesize precise yet diverse samples for various teacher models, despite minimizing \mathcal{L}_G in Eq. (2). In particular, we focus on the catastrophic failure of DAFL (Chen et al. 2019) in Figure 1(a), which heavily relies on class-prior and thus reveals the drawbacks of class-prior when training the generator, namely decreasing sample diversity yet allowing low-quality samples.

Revisiting Class-Prior in DFKD

With the goal of generating more accurate samples, the class-prior loss \mathcal{L}_{Cls} is usually defined as:

$$\mathcal{L}_{Cls} = \mathbb{E}_{z \sim p_z(z)} [\ell_{ce}(\theta_T(\theta_G(z)), \ \hat{y}_z)],$$

where \hat{y}_z is a one-hot vector corresponding to the class with the maximum probability in $\theta_T(\theta_G(z))$ and $\ell_{ce}(\cdot,\cdot)$ is the cross-entropy loss function. Since \mathcal{L}_{Cls} will continue to incur loss values with some extent until $\theta_T(\theta_G(z))$ becomes close to the one-hot vector, the generator θ_G will be more

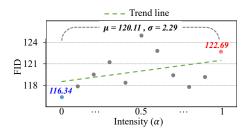


Figure 2: FID scores using a pretrained ResNet-34 model on CIFAR-10 with class-prior's intensity values from 0 to 1.



(a) With class-prior (b)

(b) Without class-prior

Figure 3: Images of *Airplane* (top) and *Dog* (bottom) generated by each trained version of the generator with or without class-prior for a pretrained ResNet-34 model on CIFAR-10.

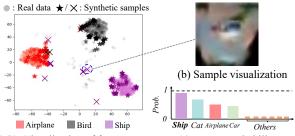
and more focused on producing less challenging samples, rather than exploring various sample cases that might be useful for transferring the teacher's knowledge. This will potentially reduce the overall diversity of generated samples, leading to less effective distillation from the teacher. Based on our intuition, this subsection conducts a detailed experimental analysis on class-prior, considering its necessity in the generator loss function.

Lower Sample Diversity. To evaluate the impact of classprior on the diversity of generated samples, we train the generator θ_G using a trained ResNet-34 model on CIFAR-10, while varying the intensity parameter of \mathcal{L}_{Cls} (i.e., α in Eq. (2)) and fixing those of \mathcal{L}_{Rep} and \mathcal{L}_{Adv} . To measure the sample diversity, we compute the Frechet Inception distance (FID) score over the samples generated by each trained version of θ_G , where the FID score (Heusel et al. 2017) is known to be smaller when evaluating more diverse and realistic samples in generative models. Figure 2 shows a roughly decreasing trend of the FID score when reducing the classprior's intensity, implying that the stronger the class-prior loss, the lower the diversity of generated samples. In Figure 3, we also visually demonstrate that a generator trained with class-prior produces a limited variety of images for Airplane and Dog classes, whereas it becomes able to generate variants of those images when removing class-prior from the generator loss function. Finally, as shown in Table 2, this trend turns out to remain the same even when using various teacher models of similar performance. With the exception of T4, where class-prior appears to be effective, the FID scores without class-prior are mostly smaller (and thus exhibit higher diversity) than those with class-prior.

Incomplete Quality Control. We next investigate how effectively the class-prior loss controls the quality of generated samples, by training a generator with class-prior but

Teachers	T1 (95.5)	T2 (93.5)	T3 (92.0)	T4 (94.4)	T5 (91.4)
$\begin{array}{c} \alpha = 0 \\ \alpha = 1 \end{array}$	116.3 122.7	119.3 122.9	106.7 111.7	109.7 105.1	106.4 109.6

Table 2: FID scores using five different ResNet-34 teacher models on CIFAR-10.



(a) 2D Visualization of feature vectors (c) Ou

(c) Output probability

Figure 4: (a) 2D visualization of feature vectors corresponding to real data and synthetic data generated by a generator trained using class-prior without the adversarial loss in ResNet-34 on CIFAR-10, where \bullet , \star , and \times represent real data samples, high-quality synthetic samples within the boundary of their corresponding real data, and low-quality ones out of their boundary. (b) and (c) show a low-quality synthetic image and its probability distribution, respectively.

removing the adversarial loss \mathcal{L}_{Adv} from Eq. (2). This is because, as pointed out by Fang et al. (2021), adversarial training seems to be the major component that causes lowquality samples in DFKD, while class-prior is supposed to enhance the sample quality. As observed in Figure 4, however, even such a generator without \mathcal{L}_{Adv} often synthesizes unexpectedly low-quality samples (represented as ×-shaped points in Figure 4(a) and visualized in Figure 4(b)) to the point that the teacher model cannot be confident about their predicted classes (see Figure 4(c)). Unfortunately, these erroneous samples have been consistently observed to account for approximately 7-8% per batch, and can confuse even well-trained teacher models, potentially leading to the failure of the entire distillation process due to their accumulated errors. Therefore, we conclude that the class-prior loss cannot solely prevent the generation of unexpectedly lowquality samples.

Proposed Method

Generator Loss Without Class-Prior. Based on the limitations of class-prior, our first remedy is to remove \mathcal{L}_{Cls} from \mathcal{L}_{G} , and therefore we have:

$$\mathcal{L}_G = \beta \mathcal{L}_{Adv} + \gamma \mathcal{L}_{Rep}. \tag{3}$$

For \mathcal{L}_{Adv} and \mathcal{L}_{Rep} , we first define their individual loss functions, $\ell_{adv}(\hat{x})$ and $\ell_{rep}(\hat{x})$, respectively, for a synthetic sample \hat{x} . Then, \mathcal{L}_{Adv} and \mathcal{L}_{Rep} are simply the expectations of their individual losses over the generated samples.

For the adversarial loss $\ell_{adv}(\hat{x})$, in common with the recent DFKD methods (Yin et al. 2020; Binici et al. 2022a;

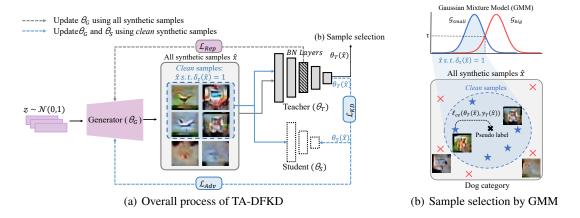


Figure 5: Overview of the proposed TA-DFKD method.

Patel, Mopuri, and Qiu 2023), we use the Jensen-Shannon (JS) divergence $JSD(\cdot,\cdot)$ as follows:

$$\ell_{adv}(\hat{x}) = 1 - JSD(\theta_T(\hat{x}), \theta_S(\hat{x})).$$

Minimizing $\ell_{adv}(\hat{x})$ maximizes the discrepancy between the outputs of the teacher and student models with \hat{x} , thereby guiding the generator to produce more difficult samples.

To specify $\ell_{rep}(\hat{x})$, we adopt the BNS technique (Yin et al. 2020), which matches the statistics of batch normalization (BN) layers to make generated samples more realistic, by the following definition:

$$\ell_{rep}(\hat{x}) = \ell_{bns}(\hat{x}) + \lambda \ell_{var}(\hat{x}) + (1 - \lambda)\ell_{l2}(\hat{x}),$$

where $\ell_{bns}(\hat{x})$ is the sum of differences between the statistics stored in BN layers of the teacher model when training real data, μ_l and σ_l^2 , and those obtained by generated samples in the teacher's same layers, $\mu_l(\hat{x})$ and $\sigma_l^2(\hat{x})$, as: $\ell_{bns}(\hat{x}) = \sum_l (\parallel \mu_l(\hat{x}) - \mu_l \parallel_2 + \parallel \sigma_l^2(\hat{x}) - \sigma_l^2 \parallel_2)$. As in the original BNS technique (Yin et al. 2020), we also leverage additional regularization terms ℓ_{var} and ℓ_{l2} , which are about total variance on pixel values within each image \hat{x} and L2-norm of \hat{x} , respectively.

By minimizing both \mathcal{L}_{Adv} and \mathcal{L}_{Rep} , the generator can effectively synthesize samples as diverse as possible and mimic the feature-level summary of real data distribution by matching BN statistics.

Quality Control via Sample Selection. Eliminating the class-prior restriction could potentially lead to an even higher risk of generating unexpectedly low-quality samples. Furthermore, the adversarial loss itself has its own problems that need to be addressed, such as the drastic change in the distribution of generated samples, as highlighted by the recent studies (Binici et al. 2022a; Patel, Mopuri, and Qiu 2023; Do et al. 2022). To address both issues, we propose a simple yet effective approach: teacher-driven sample selection, which takes only *clean* samples that are confidently verified by the given teacher model. By doing so, from any teacher models, we not only avoid distillation with erroneous samples, but also possibly mitigate drastic changes in the sample distribution.

More specifically, for a generated sample \hat{x} , we measure the quality of \hat{x} by quantifying how confident the teacher model is about its predicted label, denoted by $y_T(\hat{x})$. To this end, we compute the cross-entropy loss between the teacher's output probability and its one-hot vector of the predicted label as: $\ell_{ce}(\theta_T(\hat{x}), y_T(\hat{x}))$. This per-sample loss value is then used to determine whether \hat{x} is reliable enough in terms of its label distribution. Instead of making a decision by some absolute comparison, we specifically employ the Gaussian Mixture Model (GMM), inspired by a method of learning with noisy labels (Li, Socher, and Hoi 2020). As illustrated in Figure 5(b), for each sample batch, the GMM is built upon per-sample loss values, thereby forming two Gaussian distribution components, namely \mathcal{G}_{small} and \mathcal{G}_{big} . The \mathcal{G}_{small} component corresponds to the samples with smaller loss values, which thus are considered to be high-quality samples, while the samples belonging to \mathcal{G}_{big} are likely to be low-quality ones. To determine whether to select \hat{x} or not, we compute its posterior probability $Pr(\mathcal{G}_{small}|\ell_{ce}(\theta_T(\hat{x}), y_T(\hat{x})))$ and check if the probability exceeds a specified threshold τ . This enables us to define the following Boolean function $\delta_{\tau}(\hat{x})$:

$$\delta_{\tau}(\hat{x}) = \begin{cases} 1 & \text{if } Pr(\mathcal{G}_{small} | \ell_{ce}(\theta_{T}(\hat{x}), y_{T}(\hat{x}))) > \tau, \\ 0 & \text{otherwise.} \end{cases}$$

Finally, given a set of generated samples, we select only the subset of samples with $\delta_{\tau}(\hat{x}) = 1$ as:

$$\{\hat{x} \mid \hat{x} = \theta_G(z) \text{ s.t. } z \sim p_z(z) \land \delta_\tau(\hat{x}) = 1\}.$$

Figure 6 demonstrates the effectiveness of our sample selection in DFKD. Before applying sample selection to the generator being trained with Eq. (3), we can still observe unexpectedly low-quality synthetic samples, as indicated by x-shaped points in Figure 6(a). However, they are effectively removed from the result of Figure 6(b) after our sample selection method is applied, and therein all the synthetic samples are properly located within their corresponding boundary of real data samples. In our experiments, setting τ to 0.5, in the early stages of training, approximately 60% of samples are selected, but when approaching the end of training, more than 90% of samples are selected.

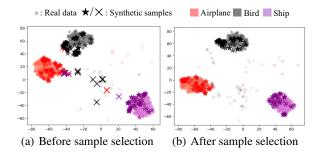


Figure 6: Visualization of before and after sample selection from real validation data and synthetic samples in the teacher model ResNet-34 on CIFAR-10.

Overall Process of TA-DFKD. We now present the overall process of our teacher-agnostic data-free knowledge distillation (TA-DFKD) method, as illustrated in Figure 5(a). When training the generator with our loss function in Eq. (3), we use all the synthetic samples without sample selection to compute \mathcal{L}_{Rep} , as \mathcal{L}_{Rep} is for learning feature-level summary of real data distribution. On the other hand, in terms of both \mathcal{L}_{Adv} and \mathcal{L}_{KD} , we train with only selected samples by our selection method. Therefore, we accordingly define the following two loss functions of Eq. (3):

$$\begin{array}{lcl} \mathcal{L}_{Adv} & = & \mathbb{E}_{z \sim p_z(z) \ \land \ \delta_\tau(\theta_G(z)) = 1}[\ell_{adv}(\theta_G(z))], \ \ \text{and} \\ \mathcal{L}_{Rep} & = & \mathbb{E}_{z \sim p_z(z)}[\ell_{rep}(\theta_G(z))]. \end{array}$$

The final KD loss is similarly defined as:

$$\mathcal{L}_{KD} = \mathbb{E}_{z \sim p_z(z) \ \land \ \delta_{\tau}(\theta_G(z)) = 1} \parallel \theta_T(\theta_G(z)) - \theta_S(\theta_G(z)) \parallel_1,$$

where we use the L1-distance between two outputs using only selected synthetic samples.

Experiments

In this section, we validate our TA-DFKD method, with a focus on its robustness and stability in DFKD using various pretrained teacher models with the similar test performance.

Environment

Datasets and Compared Methods. We use three benchmark datasets, CIFAR-10/CIFAR-100 (Krizhevsky and Hinton 2009) and TinyImageNet (Deng et al. 2009). The CI-FAR datasets contain 60,000 RGB images of 32×32 over either 10 or 100 classes, whereas Tiny-ImageNet consists of 100,000 images for 200 classes, 500 for each class, all of which is the same size of 64×64 . Using these datasets, we compare TA-DFKD with multiple SOTA DFKD methods, which are two fold. The first category includes DAFL (Chen et al. 2019) and DFAD (Fang et al. 2019), which partially mix out of the three loss terms, class-prior, adversarial, and representation losses. For the second category, we test CMI (Fang et al. 2021) as a representative one using the BNS loss, and PRE-DFKD (Binici et al. 2022a) that alternatively uses activation maximization for the same purpose. Except for CMI (Fang et al. 2021), which requires an additional training phase with pre-generated samples in memory,

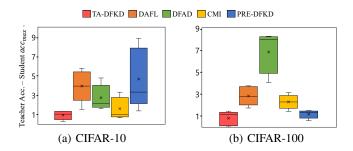


Figure 7: Box-plots of performance differences between teacher and student in the CIFAR datasets.

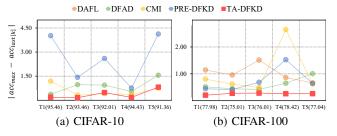


Figure 8: Differences between acc_{max} and $acc_{last[k]}$ in the CIFAR datasets.

all the compared methods are one-phase DFKD methods. Due to the space limit, the results using TinyImageNet are presented in the Appendix.

Training Details. For all datasets, we train ResNet-34 (He et al. 2016) as the teacher model, ResNet-18 as the student model, and DCGAN (Radford, Metz, and Chintala 2016) as the generator. During training ResNet-34, multiple teacher models with similar performance are randomly selected. In the entire DFKD process, we train ResNet-18 along with DCGAN for a particular number of epochs, 200 epochs for CIFAR-10 and 500 epochs for CIFAR-100 and Tiny-ImageNet. For compared methods, we follow the same configuration of their implementations. Every measurement in this section is taken out of 4 repeated runs.

Evaluation Metrics. In order to evaluate the robustness and stability of each method, we not only measure the peak accuracy acc_{max} of each student model over all the repeated runs but also introduce the converging accuracy $acc_{last[k]}$, which is the average student accuracy over the last k epochs of KD training for each run. Small differences between acc_{max} and $acc_{last[k]}$ imply that the student model shows stable and reasonably good performance during the last phase of training. Furthermore, a small deviation of $acc_{last[k]}$ out of all the repeated runs indicates a high level of the robustness within a particular teacher model. We set k to 10 for CIFAR-10 and 20 for the other datasets.

Experimental Results

Performance Comparison. Table 3 presents the summarized result of performance comparison of TA-DFKD with the SOTA DFKD methods, using five different teacher mod-

CIFAR-10 Teacher: ResNet-34 Student: ResNet-18 (Accuracy with real data: 95.2 %)										
Method	Teacher 1 ($acc_{last[10]}$	(95.46%) acc_{max}	Teacher 2 ($acc_{last[10]}$	(93.46%) acc_{max}	Teacher 3 ($acc_{last[10]}$	92.01%) acc_{max}	Teacher 4 ($acc_{last[10]}$	(94.43%) acc_{max}	Teacher 5 ($acc_{last[10]}$	91.36%) acc_{max}
DAFL DFAD CMI PRE-DFKD TA-DFKD	$ \begin{vmatrix} 83.60_{\pm 7.8} \\ \underline{93.23}_{\pm 0.1} \\ 92.54_{\pm 1.6} \\ 89.22_{\pm 4.9} \\ \textbf{94.24}_{\pm 0.1} \end{vmatrix} $	92.07 93.60 94.80 94.10 <u>94.43</u>	$ \begin{vmatrix} 85.94_{\pm 2.0} \\ 87.72_{\pm 0.2} \\ \underline{89.84}_{\pm 0.1} \\ 85.16_{\pm 0.3} \\ \textbf{91.99}_{\pm 0.1} \end{vmatrix} $	88.43 88.69 <u>90.16</u> 86.59 92.15	$ \begin{vmatrix} 68.87_{\pm 20.0} \\ 87.83_{\pm 0.1} \\ \underline{89.40}_{\pm 0.1} \\ 80.55_{\pm 0.6} \\ \textbf{90.21}_{\pm 0.1} \end{vmatrix} $	88.08 88.77 <u>89.81</u> 83.15 90.69	$\begin{array}{c} 89.21_{\pm 4.6} \\ 92.27_{\pm 0.1} \\ \underline{93.26}_{\pm 0.1} \\ \underline{90.80}_{\pm 0.3} \\ 93.61_{\pm 0.0} \end{array}$	92.90 92.82 <u>93.61</u> 91.56 93.79	$ \begin{array}{c c} 72.85_{\pm 10.9} \\ 87.66_{\pm 0.2} \\ \underline{89.62}_{\pm 0.1} \\ 83.93_{\pm 1.6} \\ \textbf{90.27}_{\pm 0.1} \end{array} $	85.59 89.22 <u>90.37</u> 88.05 91.08
	CIF	AR-100	Teacher: ResN	Net-34 Stu	ıdent: ResNet-	18 (Accura	cy with real d	ata: 77.1 %)	
Method	Teacher 1 ($acc_{last[20]}$	(77.98%) acc_{max}	Teacher 2 ($acc_{last[20]}$	(75.01%) acc_{max}	Teacher 3 ($acc_{last[20]}$	$76.01\%) \\ acc_{max}$	Teacher 4 ($acc_{last[20]}$	(78.42%) acc_{max}	Teacher 5 ($acc_{last[20]}$	$77.04\%) \\ acc_{max}$
DAFL DFAD CMI PRE-DFKD TA-DFKD	$ \begin{array}{ c c c } \hline 74.08 \pm 0.6 \\ 69.51 \pm 0.3 \\ 74.10 \pm 0.2 \\ \hline 76.13 \pm 0.2 \\ \hline 76.55 \pm 0.1 \\ \hline \end{array}$	75.22 70.03 74.85 <u>76.57</u> 76.76	$ \begin{array}{ c c c c }\hline 70.31_{\pm0.4}\\ 66.31_{\pm0.1}\\ 71.81_{\pm0.1}\\ \hline 73.11_{\pm0.2}\\ \hline \textbf{73.61}_{\pm0.1}\\ \end{array}$	71.27 66.75 72.43 <u>73.53</u> 73.89	$\begin{array}{c c} 72.22_{\pm 1.0} \\ 71.55_{\pm 0.2} \\ 73.55_{\pm 0.1} \\ \underline{74.75}_{\pm 0.4} \\ \hline \textbf{75.74}_{\pm 0.1} \end{array}$	73.73 71.97 74.03 <u>75.44</u> 76.02	$\begin{array}{c} 73.96_{\pm 0.5} \\ 69.61_{\pm 0.3} \\ 74.40_{\pm 0.1} \\ \underline{75.58}_{\pm 1.1} \\ \overline{\textbf{76.73}}_{\pm 0.1} \end{array}$	74.82 70.26 <u>77.00</u> 77.10 76.99	$\begin{array}{c} 74.66_{\pm0.4} \\ 70.33_{\pm0.4} \\ 74.18_{\pm0.1} \\ \underline{75.37}_{\pm0.6} \\ \hline \textbf{76.58}_{\pm0.1} \end{array}$	75.32 71.33 74.77 <u>76.01</u> 76.84

Table 3: DFKD performance comparison using 5 teacher models trained on CIFAR-10 (top) and CIFAR-100 (bottom).

els trained on the CIFAR-10 and CIFAR-100 datasets. It is clearly observed that TA-DFKD manages to achieve the highest peak accuracy acc_{max} as well as the highest converging accuracy $acc_{last[k]}$ in most of the cases. Over all the repeated runs, TA-DFKD shows only small variations in its converging accuracy, implying high robustness within each teacher model. On the other hand, DAFL sometimes experiences a catastrophic failure of distillation with a large deviation even with the same teacher model (e.g., \pm 20.02 in Teacher 3 on CIFAR-10), aligning with the example of Figure 1(a). This failure is likely to happen when the generator gets collapsed into only a few easy samples at some point of training. Recent methods utilizing all the three loss terms, CMI and PRE-DFKD, generally perform better than those of not using all the terms, DAFL and DFAD. However, both CMI and PRE-DFKD do not show the reliable performance across the two datasets in that either of them interchangeably takes the second best position in different datasets.

Teacher-Agnostic Behavior. Based on the results of Table 3, we examine how robust and stable the performance of each method remains when using different teacher models, as plotted in Figures 7 and 8. Figure 7 shows boxplots on performance gaps in peak accuracies between the teacher and student models. In both CIFAR-10 and CIFAR-100, the proposed TA-DFKD shows short box-plots implying the teacher-agnostic robustness, while the other compared methods have relatively long ranges of performance gaps throughout different teacher models. Figure 8 demonstrates the teacher-agnostic stability by plotting differences between acc_{max} and $acc_{last[k]}$ using five teacher models. TA-DFKD clearly takes the bottom-most position in both graphs of the CIFAR datasets, meaning that its performance becomes quite stable and remains almost the same as its best accuracy once it reaches the last phase of distillation process.

Ablation Study. Table 4 shows the result of an ablation study to verify the effectiveness of elimination of class-prior

Method	T1 (77.98)	T2 (75.01)	T3 (76.01)	T4 (78.42)	T5 (77.04)
Baseline w/o \mathcal{L}_{Cls} TA-DFKD	73.67	71.09	74.22	73.14	76.08
	73.96	<u>71.23</u>	74.43	76.08	<u>76.67</u>
	76.76	73.89	76.02	76.99	76.84

Table 4: Ablation study showing peak accuracies on CIFAR-100, where (1) baseline is the standard DFKD framework involving all the three loss terms, (2) w/o class-prior is the method removing class-prior from the standard framework, and (3) TA-DFKD is our final version additionally applying teacher-driven sample selection.

and applying sample selection, using the five teacher models on CIFAR-100. The baseline methods use all the three loss terms, \mathcal{L}_{Cls} , \mathcal{L}_{Rep} , and \mathcal{L}_{Adv} , without any sample selection. The result clearly confirms our two arguments: (1) classprior is better to be removed, but (2) removing class-prior is not sufficient to further improve the performance without controlling sample quality by our sample selection method.

Conclusion

This paper has conducted the first study on teacher-agnostic DFKD, with a focus on three loss terms commonly adopted in DFKD methodologies. Our findings strongly suggest that by replacing the class-prior restriction with our sample selection scheme, we can achieve enhanced quality control, thus leading us to propose the TA-DFKD method. In our experiments, TA-DFKD has demonstrated remarkable robustness and stability across various teacher models. We believe that our work offers a practical solution for knowledge distillation scenarios without access to prior data samples, and it is our hope that this work marks the initiation of the problem of teacher-agnostic DFKD, providing a promising direction for further research in the field.

Acknowledgments

This work was supported in part by Institute of Information & communications Technology Planning & Evaluation (IITP) grants funded by the Korea government(MSIT) (No.2022-0-00448, Deep Total Recall: Continual Learning for Human-Like Recall of Artificial Neural Networks, No.RS-2022-00155915, Artificial Intelligence Convergence Innovation Human Resources Development (Inha University)), in part by the National Research Foundation of Korea (NRF) grants funded by the Korea government (MSIT) (No.2021R1F1A1060160, No.2022R1A4A3029480), and in part by INHA UNIVERSITY Research Grant.

References

- Binici, K.; Aggarwal, S.; Pham, N. T.; Leman, K.; and Mitra, T. 2022a. Robust and Resource-Efficient Data-Free Knowledge Distillation by Generative Pseudo Replay. In *Association for the Advancement of Artificial Intelligence, AAAI*, 6089–6096. AAAI Press.
- Binici, K.; Pham, N. T.; Mitra, T.; and Leman, K. 2022b. Preventing Catastrophic Forgetting and Distribution Mismatch in Knowledge Distillation via Synthetic Data. In *Winter Conference on Applications of Computer Vision, WACV*, 3625–3633. IEEE.
- Chen, H.; Wang, Y.; Xu, C.; Yang, Z.; Liu, C.; Shi, B.; Xu, C.; Xu, C.; and Tian, Q. 2019. Data-Free Learning of Student Networks. In *International Conference on Computer Vision*, *ICCV*, 3513–3521. IEEE.
- Deng, J.; Socher, R.; Fei-Fei, L.; Dong, W.; Li, K.; and Li, L.-J. 2009. ImageNet: A large-scale hierarchical image database. In *Conference on Computer Vision and Pattern Recognition, CVPR*, volume 00, 248–255.
- Do, K.; Le, H.; Nguyen, D.; Nguyen, D.; Harikumar, H.; Tran, T.; Rana, S.; and Venkatesh, S. 2022. Momentum Adversarial Distillation: Handling Large Distribution Shifts in Data-Free Knowledge Distillation. In *Advances in Neural Information Processing Systems*, *NeurIPS*.
- Fang, G.; Song, J.; Shen, C.; Wang, X.; Chen, D.; and Song, M. 2019. Data-Free Adversarial Distillation. *CoRR*, abs/1912.11006.
- Fang, G.; Song, J.; Wang, X.; Shen, C.; Wang, X.; and Song, M. 2021. Contrastive Model Invertion for Data-Free Knolwedge Distillation. In *International Joint Conference on Artificial Intelligence, IJCAI*, 2374–2380. ijcai.org.
- Goodfellow, I. J.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A. C.; and Bengio, Y. 2014. Generative Adversarial Networks. *CoRR*, abs/1406.2661.
- Han, B.; Yao, Q.; Yu, X.; Niu, G.; Xu, M.; Hu, W.; Tsang, I. W.; and Sugiyama, M. 2018. Co-teaching: Robust training of deep neural networks with extremely noisy labels. In *Advances in Neural Information Processing Systems*, *NeurIPS*, 8536–8546.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep Residual Learning for Image Recognition. In *Conference on Computer Vision and Pattern Recognition*, CVPR, 770–778. IEEE.

- Heusel, M.; Ramsauer, H.; Unterthiner, T.; Nessler, B.; and Hochreiter, S. 2017. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. In *Advances in Neural Information Processing Systems*, *NeurIPS*, 6626–6637.
- Hinton, G. E.; Vinyals, O.; and Dean, J. 2015. Distilling the Knowledge in a Neural Network. *CoRR*, abs/1503.02531.
- Hong, J.; Zeng, Y.; Yu, S.; Lyu, L.; Jia, R.; and Zhou, J. 2023. Revisiting Data-Free Knowledge Distillation with Poisoned Teachers. In *International Conference on Machine Learning,ICML*, volume 202 of *Proceedings of Machine Learning Research*, 13199–13212. PMLR.
- Jiang, L.; Zhou, Z.; Leung, T.; Li, L.; and Fei-Fei, L. 2018. MentorNet: Learning Data-Driven Curriculum for Very Deep Neural Networks on Corrupted Labels. In *International Conference on Machine Learning,ICML*, volume 80, 2309–2318. PMLR.
- Krizhevsky, A.; and Hinton, G. 2009. Learning multiple layers of features from tiny images. Technical Report 0.
- Li, J.; Socher, R.; and Hoi, S. C. H. 2020. DivideMix: Learning with Noisy Labels as Semi-supervised Learning. In *International Conference on Learning Representations, ICLR*. OpenReview.net.
- Li, J.; Zhou, S.; Li, L.; Wang, H.; Bu, J.; and Yu, Z. 2023. Dynamic data-free knowledge distillation by easy-to-hard learning strategy. *Inf. Sci.*, 642: 119202.
- Malach, E.; and Shalev-Shwartz, S. 2017. Decoupling "when to update" from "how to update". In *Advances in Neural Information Processing Systems*, *NeurIPS*, 960–970. Micaelli, P.; and Storkey, A. J. 2019. Zero-shot Knowledge Transfer via Adversarial Belief Matching. In Wallach, H. M.; Larochelle, H.; Beygelzimer, A.; d'Alché-Buc, F.; Fox, E. B.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, *NeurIPS*, 9547–9557.
- Nayak, G. K.; Mopuri, K. R.; Shaj, V.; Radhakrishnan, V. B.; and Chakraborty, A. 2019. Zero-Shot Knowledge Distillation in Deep Networks. In *International Conference on Machine Learning, ICML*, volume 97, 4743–4751. PMLR.
- Patel, G.; Mopuri, K. R.; and Qiu, Q. 2023. Learning to Retain while Acquiring: Combating Distribution-Shift in Adversarial Data-Free Knowledge Distillation. *CoRR*, abs/2302.14290.
- Qian, B.; Wang, Y.; Hong, R.; and Wang, M. 2023. Adaptive Data-Free Quantization. In *Conference on Computer Vision and Pattern Recognition*, CVPR, 7960–7968. IEEE.
- Radford, A.; Metz, L.; and Chintala, S. 2016. Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks. In *International Conference on Learning Representations, ICLR*.
- Song, H.; Kim, M.; Park, D.; and Lee, J. 2020. Learning from Noisy Labels with Deep Neural Networks: A Survey. *CoRR*, abs/2007.08199.
- Yin, H.; Molchanov, P.; Alvarez, J. M.; Li, Z.; Mallya, A.; Hoiem, D.; Jha, N. K.; and Kautz, J. 2020. Dreaming to Distill: Data-Free Knowledge Transfer via DeepInversion. In *Conference on Computer Vision and Pattern Recognition, CVPR*, 8712–8721. IEEE.

Yoo, J.; Cho, M.; Kim, T.; and Kang, U. 2019. Knowledge Extraction with No Observable Data. In *Advances in Neural Information Processing Systems*, NeurIPS, 2701–2710.

Yu, X.; Han, B.; Yao, J.; Niu, G.; Tsang, I. W.; and Sugiyama, M. 2019. How does Disagreement Help Generalization against Label Corruption? In *International Conference on Machine Learning,ICML*, volume 97, 7164–7173. PMLR.

Appendix of "Teacher as a Lenient Expert: Teacher-Agnostic Data-Free Knowledge Distillation"

In this appendix, we first (1) present additional experimental performance comparison using 2 teacher models trained on Tiny-ImageNet, and (2) more detailed experimental results on the ablation study, and then (3) provide detailed values of Figures 7 and 8. Finally, we (4) describe all the implementation details and hyperparameter values.

Details of Experimental Results

Results of Tiny-ImageNet.

Table A1 shows the performance summary of the experiments using Tiny-ImageNet. As mentioned in (Binici et al. 2022a), we have also failed to find proper hyperparameters for CMI and DFAD, and therefore we compare TA-DFKD with only DAFL and PRE-DFKD. Similar to the results of the CIFAR datasets, our TA-DFKD method outperforms the compared methods in both the robustness and stability across two teacher models. Notably, PRE-DFKD exhibits significant performance variations in the same teacher model (e.g., \pm 9.78 in Teacher 1). This is another evidence that using both class-prior and adversarial losses may not always reach to a proper balance between diversity and sample quality.

Tiny-ImageNet Teacher: ResNet-34 Student: ResNet-18 (Accuracy with real data: 64.9 %)								
Method	Teacher 1 ($acc_{last[20]}$	acc_{max}	Teacher 2 (74.92%) $acc_{last[20]} acc_{max}$					
DAFL PRE-DFKD TA-DFKD (ours)	$\begin{array}{c c} 47.76 \pm 2.06 \\ 46.45 \pm 9.78 \\ \textbf{53.00} \pm 1.57 \end{array}$	51.33 53.16 54.84	$\begin{array}{ c c c c c c }\hline 50.60 \pm 1.84 \\ 46.13 \pm 7.50 \\ \hline 53.55 \pm 0.49 \\ \hline \end{array}$	53.20 53.68 54.52				

Table A1: Results of Tiny-ImageNet in two teacher models

Detailed Ablation Study.

We conduct an ablation study using the CIFAR-10 dataset and report the averaged outcomes from 4 repeated runs, rather than focusing only on maximum values. Similar to the results of Table 4 using the CIFAR-100 dataset, we can demonstrate the effectiveness of applying sample selection and removing class-prior.

Method		CIFAR-10								
	Teacher 1 (95.46%)	Teacher 2 (93.46%)	Teacher 3 (92.01%)	Teacher 4 (94.43%)	Teacher 5 (91.36%)					
Baseline w/o class-prior TA-DFKD (ours)	$\begin{array}{c c} 93.25_{\pm 0.14} \\ 93.06_{\pm 0.09} \\ 94.35_{\pm 0.06} \end{array}$	$90.54_{\pm 0.18}$ $90.61_{\pm 0.12}$ $92.10_{\pm 0.06}$	$88.52_{\pm 0.17}$ $88.71_{\pm 0.12}$ $90.46_{\pm 0.15}$	$\begin{array}{c} \underline{92.88}_{\pm 0.07} \\ \underline{92.84}_{\pm 0.14} \\ \underline{93.72}_{\pm 0.07} \end{array}$	$88.90_{\pm 0.15}$ $\underline{90.55}_{\pm 0.10}$ $90.85_{\pm 0.20}$					
Method			CIFAR-100							
Wethod	Teacher 1 (77.98%)	Teacher 2 (75.01%)	Teacher 3 (76.01%)	Teacher 4 (78.42%)	Teacher 5 (77.04%)					
Baseline w/o class-prior TA-DFKD (ours)	$73.67_{\pm 0.34}$ $73.96_{\pm 0.19}$ $76.74_{\pm 0.02}$	$71.09_{\pm 0.22}$ $71.23_{\pm 0.21}$ $73.83_{\pm 0.06}$	$74.22_{\pm 0.10}$ $74.43_{\pm 0.16}$ $75.93_{\pm 0.06}$	$73.14_{\pm 0.25}$ $76.08_{\pm 0.26}$ $76.32_{\pm 0.40}$	$73.64_{\pm 0.15}$ $76.67_{\pm 0.08}$ $76.79_{\pm 0.03}$					

Table A2: Detailed results of the ablation study using CIFAR-10 and CIFAR-100

Detailed Values in Figure 7 and Figure 8

Table A3 and Table A4 present detailed values corresponding to the graphs shown in Figures 7 and 8, respectively. Table A3 specifically illustrate the differences between the accuracy of the teacher and the peak accuracy of the student denoted as acc_{max} . TA-DFKD achieves the best teacher-agnostic robustness by almost always exhibiting the smallest variance from the corresponding teacher's accuracy. Table A4 shows differences between the peak accuracy (acc_{max}) and the converging accuracy $(acc_{last[k]})$ within each student model, in order to assess the stability of the student model during its converging phase of training. In most of the cases, our proposed TA-DFKD method demonstrates superior stability in performance, evidenced by minimal deviations between the best accuracy and the average accuracy over the last epochs.

Method	CIFAR-10							CIFAR-100		
	T1(95.46%)	T2(93.46%)	T3(92.01%)	T4(94.43%)	T5(91.36%)	T1(77.98%)	T2(75.01%)	T3(76.01%)	T4(78.42%)	T5(77.04%)
DAFL	3.39	5.03	3.93	1.53	5.77	2.76	3.74	2.28	3.60	1.72
DFAD	1.86	4.77	3.24	1.61	2.14	7.95	8.26	4.04	8.16	5.71
CMI	0.66	3.30	2.20	0.82	0.99	3.13	2.58	1.98	1.42	2.27
PRE-DFKD	1.36	6.87	8.86	2.87	3.31	<u>1.41</u>	1.48	0.57	1.32	1.03
TA-DFKD	1.03	1.31	1.32	0.64	0.28	1.22	1.12	-0.01	1.43	0.20

Table A3: Differences in performance between teachers and student models (Teacher Acc. - Student Peak Acc. (acc_{max})).

Method	CIFAR-10							CIFAR-100		
	T1(95.46%)	T2(93.46%)	T3(92.01%)	T4(94.43%)	T5(91.36%)	T1(77.98%)	T2(75.01%)	T3(76.01%)	T4(78.42%)	T5(77.04%)
DAFL	8.47	2.49	19.21	3.69	12.74	1.14	0.96	1.51	0.86	0.66
DFAD	0.37	0.97	0.94	0.55	1.56	0.52	0.44	0.42	0.65	1.00
CMI	2.26	0.32	0.41	0.35	0.75	0.75	0.62	0.48	2.60	0.59
PRE-DFKD	4.88	1.43	2.60	0.76	4.12	0.44	0.42	0.69	1.52	0.64
TA-DFKD	0.19	0.16	0.48	0.18	0.82	0.21	0.28	0.28	0.26	0.26

Table A4: Differences between the peak accuracy and converging accuracy within each student model (Peak Acc. (acc_{max}) - Converging Acc. $(acc_{last[k]})$.

Implementation Details

How to Train Various Teacher Models

Let us provide how we get various teacher models in detail. For the teacher models, except for Tiny-ImageNet, we use the same training environment as used in DAFL (Chen et al. 2019), referring to its implementation page: https://github.com/autogyro/DAFL/blob/master/teacher-train.py. Basically, during a relatively longer period of training, we randomly take multiple trained models as long as their performance is reasonably high.

More specifically, in MNIST, during the training of 100 epochs, we first selected a reference teacher model, so-called Teacher 1 (98.9 %), whose accuracy is the same as they are in the existing DFKD methods (Fang et al. 2019; Binici et al. 2022a). Then, we have picked up a slightly better version as Teacher 2 (99.2%).

In the CIFAR datasets, we increased the number of epochs to be 500 and adjusted the learning rate step decay by the same ratio (reducing by $0.1 \times$ at epochs 200 and 300). As in MNIST, we first selected the models with accuracy 95.46 and 77.46 as the reference teachers (i.e., Teacher 1) for CIFAR-10 and CIFAR-100, respectively, as those accuracies similarly appeared in the existing DFKD methods. For the other teacher models, we randomly selected ones whose accuracies are not below 91% for CIFAR-10 and 75% for CIFAR-100.

For Tiny-ImageNet, we fine-tuned a pretrained model on ImageNet with the Tiny-ImageNet dataset to the point that the resulting accuracy of the model gets reasonably high, such as 75%, and randomly selected two teacher models.

Details of Implementation and Hyperparameters of TA-DFKD

Our code is available at the following website: https://github.com/bigdata-inha/TA-DFKD-Official. For fair comparison, we follow the same settings, which include the number of epochs and iterations per epoch for the entire DFKD process, as those of PRE-DFKD (Binici et al. 2022a). For all datasets, we use the same set of hyperparameters. More specifically, we set the dimensionality of latent vectors (i.e., z) to 1000, and set the batch size to 1024. When training the generator, we adopted Adam optimizer with $\beta = 1$, $\gamma = 10$ and a learning rate of 0.001. In the distillation of student models, based on SGD optimizer, we used 0.01 initial learning rate with a cosine learning rate decay, 5e-4 weight decay and 0.9 momentum.