De-confounded Data-free Knowledge Distillation for Handling Distribution Shifts

Yuzheng Wang^{1*} Dingkang Yang^{1*} Zhaoyu Chen¹ Yang Liu¹ Siao Liu¹ Wenqiang Zhang² Lihua Zhang^{1†} Lizhe Qi^{1,2,3†}

¹Shanghai Engineering Research Center of AI & Robotics, Academy for Engineering & Technology, Fudan University

²Engineering Research Center of AI & Robotics, Ministry of Education, Academy for Engineering & Technology, Fudan University ³Green Ecological Smart Technology School-Enterprise Joint Research Center

{yzwang20, dkyang20}@fudan.edu.cn

Abstract

Data-Free Knowledge Distillation (DFKD) is a promising task to train high-performance small models to enhance actual deployment without relying on the original training data. Existing methods commonly avoid relying on private data by utilizing synthetic or sampled data. However, a long-overlooked issue is that the severe distribution shifts between their substitution and original data, which manifests as huge differences in the quality of images and class proportions. The harmful shifts are essentially the confounder that significantly causes performance bottlenecks. To tackle the issue, this paper proposes a novel perspective with causal inference to disentangle the student models from the impact of such shifts. By designing a customized causal graph, we first reveal the causalities among the variables in the DFKD task. Subsequently, we propose a Knowledge Distillation Causal Intervention (KDCI) framework based on the backdoor adjustment to de-confound the confounder. KDCI can be flexibly combined with most existing state-of-the-art baselines. Experiments in combination with six representative DFKD methods demonstrate the effectiveness of our KDCI, which can obviously help existing methods under almost all settings, e.g., improving the baseline by up to 15.54% accuracy on the CIFAR-100 dataset.

1. Introduction

Deep Neural Networks (DNNs), as a powerful and reliable tool, are increasingly expected to be applied to practical artificial intelligence scenes [1–7]. Despite significant progress, good performance of deep learning models is often inseparable from large-scale models [8–13] and highquality original training data [14–20]. The dependencies hinder the deployment of this technology on mobile devices



Figure 1. Diagrams of the distribution shifts between the original and substitute data for existing DFKD methods on CIFAR-10. (a) represents the random visualization and FID score of the synthetic data by DAFL, DeepInv, and sampled by DFND. (b) indicates the proportion of sample numbers in various classes (%) of the original and substitute data.

and data privacy scenes. Therefore, model compression and data-free technology have become the key to breaking through the bottleneck. To this end, Lopes *et al.* [21] propose the Data-Free Knowledge Distillation (DFKD) task. In this process, knowledge is transferred from the cumbersome model to a small model that is more suitable for deployment [22–24] without relying on the original training data. As a result, DFKD has received more attention due to its convenience and wide application.

Since the original training data is not available for privacy or other reasons [25], the key is how to supplement the new training data, *i.e.*, the substitution data. Based on the source of the substitution data, almost all existing DFKD methods can be divided into generation-based and sampling-based methods. Despite the impressive improvements achieved by these DFKD methods through complex loss stacking [26, 27] and knowledge distillation strategies [28, 29], the trained students still suffer from distribution shifts between the substitution and original data, which has long been overlooked. First, the quality of the synthetic or sampled images significantly differs from the original. Besides, for generation-based methods, the synthetic data

^{*}Equal contribution. [†]Corresponding authors.

relies on the teacher's guidance, and it is easier to synthesize the class familiar to the generator. For samplingbased methods, the sampled data entirely depends on the teacher's preference for various classes. These protocols make the preference of the teacher model inevitably affect class proportions and also lead to distribution shifts. Such shifts confound the student learning process. For example, if a pre-trained teacher model is not familiar with a specific class A, *i.e.*, it is difficult to obtain high confidence, resulting in fewer synthetic or sampled data belonging to A. For the class balance, the teacher tends to classify ambiguous and indistinguishable data into A, leading to the distribution shifts [30]. Relying on these data, the student is inevitably confused with the original testing data with the different distributions.

More intrigued, we select three DFKD methods (DAFL [31], DeepInv [26], and DFND [32]) and perform a toy experiment on the CIFAR-10 [33]. This toy experiment aims to show the distribution shifts between the substitution and original data. These methods include generation with generators (DAFL), generation through teacher model inversion (DeepInv), and sampling based on teacher preferences (DFND). We use the original data as a comparison benchmark and compare them from two aspects: the quality of images and class proportions. The results are shown in Figure 1. In Figure 1a, we randomly visualize the original data, the substitution data of DAFL, DeepInv, and DFND, and calculate the Fréchet Inception Distance (FID, lower is better) [34], a metric widely used to evaluate the quality of images. The substitution and original data are different for the data distribution domain. In Figure 1b, we test the class proportions (the substitution data is based on teacher pseudo-labels). A prominent result is that the classes of the substitution data are unbalanced due to teacher preferences, which greatly differ from the original data. These observations confirm the distribution shifts between the substitution and original data, confounding the student model.

Based on these observations, we attempt to introduce a new perspective with causal inference to handle the distribution shifts. During the application of theoretical causal inference [35] to the DFKD task, the challenges lie in describing and designing plausible causal effects and identifying and compensating for biased student learning on the substitution data with shifts. To this end, this paper attempts to address the challenges by drawing on instinctive human causalities [36] to find causal relationships among the variables in the DFKD task and optimize the biased student training process. We first disentangle the causalities and customize the causal graph according to the properties of the variables in the DFKD task. Based on this, we explore the causal paths from the substitution inputs X to the student predictions S. Then, we propose a simple yet effective Knowledge Distillation Causal Intervention (KDCI)

framework to achieve de-confounded DFKD and use the do-calculus P(S|do(X)) to calculate the actual causal effect, instead of classic likelihood P(S|X) without considering the shifts. KDCI can be easily combined with existing methods and use the backdoor adjustment [37] to de-confound and alleviate the impact of the shifts. Experiments on KDCI combined with six representative DFKD methods demonstrate its strong positive effect on the existing DFKD pipeline. Specifically, the primary contributions and experiments are summarized below:

- To our best knowledge, we are the first to alleviate the dilemma of the distribution shifts in the DFKD task from a causality-based perspective. Such shifts are regarded as the harmful confounder, which leads the student to learn misleading knowledge.
- We propose a KDCI framework to restrain the detrimental effect caused by the confounder and attempt to achieve the de-confounded distillation process. Besides, KDCI can be easily and flexibly combined with existing generation-based or sampling-based DFKD paradigms.
- Extensive experiments on the combination with six DFKD methods show that our KDCI can bring consistent and significant improvements to existing state-of-theart models. Particularly, it improves the accuracy of the DeepInv [26] by up to 15.54% on the CIFAR-100 dataset.

2. Related Work

Data-Free Knowledge Distillation. Data-free knowledge distillation is a promising task to train small models while avoiding leakage of original training data [21, 38]. The critical point is how to supplement substitution data [39-42]. The existing methods are mainly divided into three types: Generative Adversarial Networks (GANs) generation [29, 31], teacher-based model inversion generation [26, 27], and unlabeled data sampling [28, 32, 43]. Chen et al. [31] introduce the generator into the DFKD task and improve teachers' familiarity with generating data. Fang et al. [29] propose feature sharing to simplify the generation process. To better generation quality, Yin et al. [26] explore the prior knowledge of the data. Fang et al. [27] introduce contrastive learning to enhance student performance. Chen et al. [32] and Fang et al. [28] select wild data and outof-domain (OOD) data to reduce generation costs. Despite the promising performance, a long-overlooked issue is data distribution shifts, *i.e.*, the distribution bias of the student's training data and the original data is a confounder that significantly causes performance bottlenecks.

Causal Inference. Causal inference is a theory-oriented tool that seeks actual effects in a specific phenomenon [35], which has been studied and followed by diverse fields such as economics [44] and psychology [45] communities. The mainstream causal inference studies applied to neural information processing consist of two aspects: intervention



Figure 2. The causal graph. (a) The existing methods ignore distribution shifts. (b) The shifts are alleviated by causal inference.

[46–51] and counterfactuals [52–55]. Intervention is a technique for manipulating the original data distribution to reveal causal effects [37]. Counterfactual describes the imagined results generated by factual variables when treated differently [56]. Benefiting from the strong potential of causal inference to decouple spurious correlations among variables, it is gradually adopted to improve the performance of models for different downstream tasks, such as visual question answering [57], emotion recognition [58], and scene graph generation [52]. In contrast, to our best knowledge, this is the first work to identify the distribution shifts in the DFKD task through the causal intervention and alleviate the confounding effect caused by the shifts.

3. Methodology

3.1. Causal Graph of DFKD Task

First, we customize the causal graph according to the properties of the variables in the DFKD task. Specifically, the teacher is pre-trained with original training data, which is not disturbed by distribution shifts. For the student, it uses the substitution data to train while testing on the original data. The data distribution shifts indicate that it will be disturbed by the biased data [30]. During the distillation process, the teacher and student are fed the same substitution data. In this case, the student's predictions are constrained to learn the teacher's predictions. Following the same graphical notation as [59] for clarity and interpretability, we denote the variables with the notes \mathcal{N} and construct the direct causal effects with the links \mathcal{E} . From Figure 2, there are four variables involved in the DFKD causal graph $\mathcal{G} = \{\mathcal{N}, \mathcal{E}\}$, which includes the substitution inputs X, the confounder Z, the teacher's predictions T, and the student's predictions S. In particular, our causal graph is applicable to almost all existing DFKD methods so that it can be used as a general framework. The details of the causal relationships are described as follows.

 $Z \rightarrow X$. Existing DFKD methods rely on teacher predictions to supplement substitution data. For the generation-based methods, the generator is guided by the teacher and more inclined to synthesize data that is easier to synthesize [26, 27, 29, 31]. For the sampling-based methods, the data that the teacher is most [32] or least [28] fa-

miliar with is sampled. On the one hand, these synthetic or sampled data are always class-imbalanced. On the other hand, these sources of substitution data rely heavily on the teacher, so they are highly volatile and vulnerable to teacher preferences. These issues cause the distribution shifts between the original and substitution data. The shifts are treated as the harmful confounder Z [56]. On this basis, the confounder Z causes the substitution data X to be biased compared to the original data, *i.e.*, $Z \to X$.

 $Z \rightarrow S$. Due to the distribution shifts between the substitution and original data, the student trained on the substitution data tends to produce and exhibit biased predictions during the testing stage. The detrimental confounder Z confounds and affects the student's training via the causal link $Z \rightarrow S$, which causes the performance bottleneck.

 $X \to T/S \& T \leftrightarrow S$. As with existing DFKD methods, both teacher and student make predictions on the substitution data X simultaneously. By constraining their prediction distributions, the student's parameters are updated for optimization. In our DFKD causal graph, the prediction processes of the teacher and student are represented as $X \to T$ and $X \to S$. The link $T \leftrightarrow S$ reflects the interaction causal effect between these two predictions during knowledge distillation. Through these paths, the student can learn consistent knowledge from its teacher.

According to the causal theory [35], the confounder Z as a common cause directly or indirectly impacts the substitution inputs X and the student's predictions S simultaneously. The knowledge transfer process from T to S increases the student's familiarity with these substitution data. However, the confounder Z causes X to shift the original data distribution, leading to impure knowledge, which adversely affects student performance. The detrimental effects follow the backdoor causal path as $X \leftarrow Z \rightarrow S$.

3.2. Causal Intervention via Backdoor Adjustment

In the existing DFKD task, the pre-trained teacher model is fixed while the student model is learnable. As shown in Figure 2a, existing methods rely on the likelihood estimation of the student model as P(S|X). The knowledge transfer process is expressed as:

$$P(\boldsymbol{S}|\boldsymbol{X}) = \sum_{\boldsymbol{z}} P(\boldsymbol{S}|\boldsymbol{X}, KD\langle \boldsymbol{T} = f_T(\boldsymbol{X}), \boldsymbol{S} = f_S(\boldsymbol{X}, \boldsymbol{z})\rangle)P(\boldsymbol{z}|\boldsymbol{X}),$$
(1)

where $KD\langle , \rangle$ represents the knowledge distillation process between T and S. $f_T(\cdot)$ and $f_S(\cdot)$ represent the teacher model and the student model. The confounder Z introduces the data distribution shifts via P(z|X), which makes the knowledge learned by the student impure. To get rid of the confounding effect caused by Z, an intuitive idea is changing inputs X to overcome the data distribution shifts and make X unaffected by Z, *i.e.*, we have to use the data from the same distribution with the original training set as the student's training data. However, it is not pos-



(a) Confounder dictionary construction

(b) Knowledge distillation with bias compensation

Figure 3. The overview of our KDCI. In stage (a), all substitution data is fed a pre-trained model to explore the prior knowledge and construct the confounder dictionary. In stage (b), the prototype integration is built by the confounder dictionary and is used to compensate for biased student predictions. The distillation loss is calculated between the teacher's prediction and the student's compensated prediction.

sible under the setting of the DFKD task. To tackle this issue, we introduce the backdoor adjustment [35] to construct causal intervention P(S|do(X)) and block the backdoor path between X and S via Z. As a theoretical operation, implementing backdoor adjustment can be viewed as measuring the distribution shifts by estimating the average causal effect based on the class proportions. By compensating for shifted student predictions, we alleviate the shift issue and suppress the disturbance of the confounder Z. In this case, the causal path from Z to X is cut-off in Figure 2b. The student learns pure knowledge with causal intervention P(S|do(X)) rather than original biased likelihood P(S|X). This process can be expressed as:

$$P(\boldsymbol{S}|do(\boldsymbol{X})) = \sum_{\boldsymbol{z}} P(\boldsymbol{S}|\boldsymbol{X}, KD\langle \boldsymbol{T} = f_T(\boldsymbol{X}), \boldsymbol{S} = f_S(\boldsymbol{X}, \boldsymbol{z})\rangle)P(\boldsymbol{z}),$$
(2)

where X is no longer disturbed by z since causal intervention forces X to integrate each z fairly into the predictions of S, according to the corresponding prior P(z).

3.3. De-confounded DFKD with KDCI

To de-confound the DFKD task, we propose a Knowledge Distillation Causal Intervention (KDCI) framework to alleviate the distribution shift issue. The overview of KDCI is shown in Figure 3, which contains two stages: *confounder dictionary construction* and *knowledge distillation with bias compensation*. First, after obtaining the substitution data and before training the student, we model the prior knowledge of these substitution data through the prototype clustering algorithm to obtain an intervention-driven confounder dictionary. Then, the biased student predictions are compensated based on the subcenters and proportions. Notably, for a general DFKD pipeline, our framework can be easily combined with other methods. The implementation of KDCI is as follows.

Confounder Dictionary Construction. Since the substi-

tution data has no ground-truth information and the actual classes are ambiguous, we define a confounder dictionary $Z = [z_1, z_2, \dots, z_N]$ to explore the prior knowledge of these data. N is a hyperparameter representing the confounder size and $z_i \in \mathbb{R}^d$ is a single prototype. The prior knowledge implies the potential shifts and the differentiation information of class proportions. From Figure 3a, all substitution data is fed to an experienced pre-trained model (e.g., the teacher model itself) to obtain the prediction feature set $M = \{m_j \in \mathbb{R}^d\}_{j=1}^{N_m}$, where N_m is the number of the substitution data. We employ the K-Means++ with principle component analysis as the prototype clustering algorithm. After clustering, each z_i represents a prototype feature cluster, and the prototype subcenter is put into the confounder dictionary as a prototype representation. The feature cluster is denoted as $\sum_{k=1}^{N_i} m_k^i$ and the subcenter is denoted as $z_i = \frac{1}{N_i} \sum_{k=1}^{N_i} m_k^i$, where N_i is the number of the prediction features in *i*-th cluster. Therefore, the prototype proportion can be calculated as $P(z_i) = N_i/N_m$.

Knowledge Distillation with Bias Compensation. After confounder dictionary construction, we approximate a theoretical causal inference by the confounder dictionary and prototype proportions to compensate for biased student predictions to learn pure knowledge, as shown in Figure 3b. In practice, the calculation of P(S|do(X)) requires multiple forward passes of all z resulting in expensive computational costs. To simplify the above process, we apply the Normalized Weighted Geometric Mean (NWGM) [60] and approximate the Eq. (2) as:

$$P(\boldsymbol{S}|do(\boldsymbol{X})) \approx P(\boldsymbol{S}|\boldsymbol{X}, KD\langle f_T(\boldsymbol{X}), \sum_{\boldsymbol{z}} f_S(\boldsymbol{X}, \boldsymbol{z})P(\boldsymbol{z})\rangle).$$
(3)

During the knowledge transfer process, the update of student model parameters depends on the difference in predictions between the teacher and student, *e.g.*, calculating the Kullback-Leibler (KL) divergence as $f_S \leftarrow \eta \nabla_s KL(T, S)$, where η denotes learning rate, and ∇_s denotes the gradient. Considering the distribution shift of training data, we introduce the prepared prior information of the cofounder dictionary to optimize the above process. Based on this, the student predictions after compensation are represented as the integration of the biased predictions and the prior information as: $P(S|do(X)) = \phi(f_S(X), F(z))$, where $\phi(\cdot)$ is a practically simple yet empirically powerful addition fusion strategy. The prior information F(z) is calculated as:

$$F(\boldsymbol{z}) = \sum_{i=1}^{N} \lambda_i \boldsymbol{z}_i P(\boldsymbol{z}_i), \qquad (4)$$

where λ_i is a weight coefficient that measures the importance of each prototype subcenter z_i . $P(z_i)$ is the proportion of data in the *i*-th cluster. Here, we design an implementation of λ_i with the additive attention as:

$$\lambda_i = softmax(\boldsymbol{W}_t \cdot Tanh(\boldsymbol{W}_q f_S(\boldsymbol{X}) + \boldsymbol{W}_k \boldsymbol{z}_i)), \quad (5)$$

where $W_t \in \mathbb{R}^{d_n \times 1}$, $W_q \in \mathbb{R}^{d_n \times d_h}$, and $W_k \in \mathbb{R}^{d_n \times d}$ are learnable mapping matrices.

4. Experiments

4.1. Datasets and Models

Datasets. We evaluate the proposed framework on widely used classification datasets: CIFAR-10 [33], CIFAR-100 [33], Tiny-ImageNet [61], and ImageNet [16]. CIFAR-10 and CIFAR-100 contain 50,000 training samples and 10,000 testing samples of 32×32 resolution. Tiny-ImageNet contains 100,000 training samples, 10,000 validating samples, and 10,000 testing samples of 64×64 resolution. ImageNet contains 1000 classes with 1.28 million training samples and 50,000 validating samples of 224×224 resolution. Models. We test the performance of various DFKD methods on several network architectures, including resnet [1], vgg [62], and wide resnet [63]. For CIFAR-10 and CIFAR-100, we use the pre-trained teacher models from CMI [27], unify the teacher models among all methods, and set up five teacher-student backbone combinations following ex-

isting settings [27–29]. For Tiny-ImageNet, we train a renset-34 teacher model without the mixup data augmentation [64]. And the student utilizes the renset-18 as its backbone. For ImageNet, we choose the same pre-trained resnet-50 teacher model with [65] for all baseline methods.

4.2. Method Zoo

To comprehensively verify the effectiveness of KDCI, we select representative DFKD methods, including generationbased and sampling-based methods. The generation-based methods spend extra computing costs to obtain substitute data by generative adversarial networks and teacher inversion, including DAFL [31], Fast [29], CMI [27], and Deep-Inv [39]. The sampling-based methods use unlabeled data as the substitute data, including Mosaick [28] and DFND [32]. For DAFL, Fast, and DeepInv, we follow the same settings as their original papers. For CMI, due to the unpublished pre-inversion data, we choose the base version of CMI, which leads to the performance slightly lower than that reported in the original paper. For Mosaick and DFND, we sample 600k unlabeled data in ImageNet [16] for CIFAR and Tiny-ImageNet, and 600k unlabeled data in Flicker1M dataset for ImageNet. Due to the image quality, the reported performance of Mosaick is slightly better than the original paper. The implementation details and loss functions of all the above methods are shown in *Supplementary Sec.7*.

4.3. Confounder Setup

We use a pre-trained model to obtain the prediction feature set M. By default, the pre-trained model is the teacher itself, which is trained on original data. Each prediction feature m is extracted from the logits output of the last layer, and the hidden dimension d is equal to the number of classes. By default, the number of clusters N is the same across different datasets. For the substitution data in a minibatch of model inversion [27, 39], the number of clusters N is 32. For the synthetic mini-batch from GANs [29, 31], the number of clusters N is 8. For the unlabeled substitution data in sampling methods [28, 32], the number of clusters N is 128. Due to different training paradigms, the way KDCI is combined with these methods is different. For the generation-based process, the generator and student models are updated alternately. We use a mini-batch of synthetic training data to construct the cofounder dictionary, and the dictionary will be updated as the generator is updated. For the sampling-based process, unlabeled data only needs to be filtered once. We build the confounder dictionary once before distillation. Pseudocode for the above processes and other training settings are shown in Supplementary Sec.1.

4.4. Performance Comparison

To verify the proposed KDCI framework, we compare the original version and their KDCI-based version.

Results on CIFAR-10 and CIFAR-100. The results in Table 1 show the following vital observations. (i) KDCI consistently improves the performance of existing methods on all baselines across two datasets. (ii) For CIFAR-10, although the original students' performance is already close to their teachers', KDCI still provides promising gains (mostly 1%-2% improvement) for students by eliminating the harmful impact of confounder. For some baselines with poor results, KDCI brings significant improvement, *e.g.*, up to $8.85\%^{\dagger}$ for DAFL. (iii) For CIFAR-100, KDCI can significantly improve various SOTA methods (about 3%-5% improvement on average). Under some settings, KDCI improves the original methods with slightly lower performance to competitive performance, *e.g.*, 15.54%[‡] and

Dataset			CIFAR-10		CIFAR-100						
T.backbone	resnet-34	vgg-11	wrn-40-2	wrn-40-2	wrn-40-2	resnet-34	vgg-11	wrn-40-2	wrn-40-2	wrn-40-2	
S.backbone	resnet-18	resnet-18	wrn-16-1	wrn-40-1	wrn-16-2	resnet-18	resnet-18	wrn-16-1	wrn-40-1	wrn-16-2	
Teacher	95.70	92.25	94.87	94.87	94.87	78.05	71.32	75.83	75.83	75.83	
Student	95.20	95.20	91.12	93.94	93.95	77.10	77.10	65.31	72.19	73.56	
DAFL	92.22	81.10	65.71 [†]	81.33	81.55	74.47	54.16	20.88 [♯]	42.83	43.70	
DAFL+ KDCI	92.62	81.31	74.56 [†]	82.91	82.65	74.51	58.79	31.75 [♯]	46.16	48.48	
Fast	94.05	90.53	89.29	92.51	92.45	74.34	67.44	54.02	63.91	65.12	
Fast+ KDCI	94.56	91.16	89.62	93.09	92.85	75.10	68.97	54.69	67.09	68.12	
CMI	94.24	91.24	89.16	91.93	92.00	74.64	66.68	55.28	63.44	64.22	
CMI+ KDCI	94.43	91.28	89.52	92.84	92.73	75.07	69.07	57.19	67.47	67.68	
DeepInv	93.26	90.36	83.04	86.85	89.72	61.32 ^は	54.13 [‡]	53.77	61.33	61.34	
DeepInv+ KDCI	93.67	91.42	83.47	89.32	91.06	74.59 ^は	69.67 [‡]	55.22	62.13	65.90	
Mosaick	95.27	91.69	90.03	93.28	92.94	75.91	71.58	59.32	66.61	67.36	
Mosaick+ KDCI	95.43	92.36	92.25	94.45	94.20	77.06	71.86	62.03	72.19	72.39	
DFND	95.36	91.86	90.26	93.33	93.11	74.42	68.97	59.02	69.39	69.85	
DFND+ KDCI	95.44	92.54	92.4 7	94.43	94.43	77.09	72.12	66.37	74.20	74.52	

Table 1. The accuracy (%) on CIFAR-10 and CIFAR-100 about baseline methods vs. their KDCI-based version. **T.backbone** and **S.backbone** represent the backbones of the teacher and student. **Teacher** and **Student** refer to scratch training on original data. The improved results are marked in **bold**. $\{\dagger, \ddagger, \ddagger, \ddagger\}$ denote the provenance mentioned in the analysis.

Table 2. The accuracy (%) on Tiny-ImageNet dataset. The teacher uses resnet-34, and the student uses resnet-18 as the backbones. The teacher achieves an accuracy of 52.74%. The GPU time indicates the training time of one epoch on a single RTX 3090 GPU.

Method	Accuracy (%)	GPU time	Memory-Usage
Fast	28.79	101.67s	5745M
Fast+ KDCI	38.23 (+9.44)	104.43s (+2.71%)	5748M (+0.05%)
DeepInv	20.68	255.26s	3312M
DeepInv+ KDCI	34.84 (+14.16)	258.51s (+1.27%)	3316M (+0.12%)
DFND	42.64	129.16s	4196M
DFND+ KDCI	49.54 (+6.90)	133.42s (+3.30%)	4198M (+0.05%)

 $13.27\%^{\ddagger}$ for DeepInv & $10.87\%^{\ddagger}$ for DAFL. These strong gains demonstrate that KDCI can compensate for biased student predictions to learn pure knowledge by constructing prior knowledge on the substitution data whose data distribution differs from the original data distribution. (iv) We notice a small increase for KDCI-based Fast & CMI. The reasonable explanation is that they extract prior knowledge about the substitution data by accessing the statistics in the teacher's Batch Normalization layers [66], which implicitly apply the likelihood estimation and weaken our causal intervention. (v) Besides, we are pleasantly surprised to find that the students trained by sampling-based methods (e.g., Mosaick & DFND) can slightly outperform the teacher in some settings (e.g., vgg-11 \rightarrow resnet-18), both the original and KDCI-based versions. Both Mosaick and DFND utilize the unlabeled data. With the additional rich semantic knowledge, more students outperform their teachers with the help of KDCI framework.

Results on Tiny-ImageNet. For the Tiny-ImageNet, we conduct experiments with Fast, DeepInv, and DFND. The results are shown in Table 2. With the help of KDCI, the accuracy of the three methods is increased by 9.44%, 14.16%,

Table 3. The accuracy (%) on ImageNet dataset. " \rightarrow " denotes the teacher's (left) and student's (right) backbone pair.

Settings	$ $ resnet-50 \rightarrow resnet-18	$ $ resnet-50 \rightarrow mobilenetv2
Fast	53.45	43.02
Fast+ KDCI	58.24 (+4.79)	50.12 (+7.10)
Deeplnv	51.36	40.25
Deeplnv+ KDCI	55.27 (+3.91)	46.24 (+5.99)
DFND	42.82	16.03
DFND+ KDCI	51.26 (+8.44)	34.32 (+18.29)

and 6.90%, respectively. The Tiny-ImageNet dataset contains richer semantic information, which helps construct more expressive confounders and facilitates KDCI to bring more sufficient gains. Besides, we test and show the additional calculation and memory overhead. The overhead introduced by KDCI mainly comes from the confounder matrix. The additional overhead can be almost negligible since only a simple clustering algorithm is used.

Results on ImageNet. For the ImageNet, we conduct two backbone combinations with three baseline methods. The results are shown in Table 3. The generation-based methods (Fast & Deeplnv) have to train 1,000 generators (one generator for one class). We speculate that a possible reason why KDCI has smaller gains for these two generation-based methods is that 'one generator for one class' may alleviate the distribution shifts issue to a certain extent and thereby weaken the effect of causal intervention. In comparison, the gain of KDCI for DFND is higher. Overall, from the experimental results of ImageNet, the positive impact of KDCI on students is also consistent. These results further validate the effectiveness of our method.

Combining the performance on the above datasets, we conclude that KDCI can provide more significant help on more complex datasets (*e.g.*, ImageNet & Tiny-ImageNet

Table 4. Ablation studies about the prior information $F(z) = \sum_{i=1}^{N} \lambda_i z_i P(z_i)$ in Eq. (4). The results include (1) original F(z), (2) random weight coefficient λ_i , (3) random confounder dictionary z_i , and (4) without (w/o) prototype proportion $P(z_i)$.

Settings	(1) Original F	$r(\boldsymbol{z})$	(2) Random	λ_i	(3) Random	z_i	(4	4) w/o $P(\boldsymbol{z_i})$)
Methods CIFAR-10 CIFAR-100	Fast 94.56 75.10	DeepInv 93.67 74.59	DFND 95.38 77.09	Fast 93.92 74.79	DeepInv 91.56 72.72	DFND 95.28 76.86	Fast 93.35 73.76	DeepInv 91.84 72.81	DFND 94.94 76.14	Fast 93.70 74.60	DeepInv 92.76 72.66	DFND 95.11 76.97
96 95 94 53		Fast	- CIF	95 AR-10 94 93		DeepInv		97 96 95 95 94		DFND	CIFAI	R-10
⁴ ⁴ ⁷⁵ ⁷⁶ ⁷⁶	8	16	32 — CIFA 32	92 R-100 75 74 73 64	8	16 16 16 N	32 — CIF 32	64 78 32 78 78 76 76 76 76 76 32 76 76 76 32 76 76 32 76 76 32 76 76 76 32 76 76 76 76 76 76 76 76 76 76 76 76 76	64	128 128 128 N	256 CIFAR 256	512 -100 512

Figure 4. The test accuracy (%) on CIFAR-10 and CIFAR-100 datasets about different confounder dictionary size N. The teacher uses resnet-34, and the student uses resnet-18 as the backbones.

with more classes and various visual effects). More complex datasets are more susceptible to teacher preferences, leading to more severe distribution shifts. Further, the detrimental shifts inevitably lead to biased substitution data compared to the original data. Fortunately, KDCI favorably de-confound the biased student predictions, achieving significant performance improvements.

4.5. Analysis of Prior Information F(z)

We conduct ablation studies to validate the effectiveness of the components of prior information F(z) in Eq. (4) used to compensate students for biased predictions in Table 4. We select three methods (Fast, DeepInv, and DFND) on both CIFAR-10 and CIFAR-100 datasets. The teacher and student use resnet-34 and resnet-18 as their backbones, respectively. Other settings are the same as Table 1.

Necessity of Weight Coefficient λ_i . The weight λ_i represents the degree of each confounder. Comparing (1) and (2), the random λ_i causes a decline in performance. Such results indicate that depicting the importance of each confounder is essential to achieve effective causal intervention.

Rationality of Confounder z_i . The confounder z_i comes from the predicted feature representation of the pre-trained model, which directly implies prior knowledge about the substitution data. Comparing (1) and (3), students using our custom confounder significantly outperform the alternative confounder that are randomly initialized, which proves the validity of extracted prior knowledge.

Impact of Prototype Proportion $P(z_i)$. The prototype proportion $P(z_i)$ denotes the frequency of each confounder containing the knowledge of feature proportions. From (1) and (4), the proportion of each confounder plays a vital role in precise intervention implementation.

4.6. Analysis of Confounder Dictionary Z

The confounder dictionary Z is proposed to explore the prior knowledge of the substitution data. We investigate the effectiveness of Z in two perspectives: the confounder prototype size N and the selected pre-trained models. For the size, we select representative methods to test the effect of different N. For the selected pre-trained models, we use the models coming from other datasets with different numbers of classes. We swap the pre-training models on CIFAR-10 and CIFAR-100 to build the confounder and align the feature dimensions through a learnable mapping matrix.

Impact of Confounder Dictionary Size N**.** To justify the size N of the confounder Z, we set five sets of N for each method. For Fast and DeepInv, Z comes from a mini-batch synthetic data. For DFND, Z comes from the sampled data. In Figure 4, designing the suitable N for methods that suffer from varying degrees of harmful shifts helps to perform deconfounded training better.

Impact of Confounder Dictionary Sources. Table 5 shows three settings with/without confounder dictionary Z. We have two interesting discoveries. (i) First, an obvious conclusion is that using Z outperforms the original DFKD methods without Z in almost all settings. Such observations demonstrate the effectiveness of causal intervention. (ii) Second, swapping the confounders from CIFAR-10 and CIFAR-100 teacher models brings the performance decrease. For CIFAR-10, the distribution of the substitution data is simple. Simple distributions are over-separated when features are extracted using pre-trained models from complex distributions. We call this phenomenon overintervention. The excessive causal intervention potentially causes the deviation of the confounder itself. For CIFAR-100, the distribution is more complex. The complex distributions are not well approximated when using pre-trained models with less discriminative ability. We call this phe-

Table 5. Ablation studies about the confounder dictionary Z. "w/o Z" denotes the vanilla version of DFKD methods. "original Z" denotes the original confounder from the teacher itself. "other Z" denotes the confounder from another pre-trained model, *i.e.*, swapping the confounder from the pre-training teacher models on CIFAR-10 and CIFAR-100 datasets.

Dataset	CIFAR-10					CIFAR-100						
Settings	resnet-34 \rightarrow resnet-18 vgg-11 \rightarrow resnet			$gg-11 \rightarrow resnet$	-18	resnet-34 \rightarrow resnet-18			$vgg-11 \rightarrow resnet-18$			
Z	w/o <i>Z</i>	original $oldsymbol{Z}$	other $oldsymbol{Z}$	w/o <i>Z</i>	original $oldsymbol{Z}$	other $oldsymbol{Z}$	w/o Z	original $oldsymbol{Z}$	other $oldsymbol{Z}$	w/o <i>Z</i>	original $oldsymbol{Z}$	other Z
Fast DeepInv DFND	94.05 93.26 95.36	94.56 93.67 95.44	93.96 93.56 95.41	90.53 90.36 91.86	91.16 91.42 92.54	90.73 91.26 92.34	74.42 61.32 74.34	75.10 74.59 77.09	74.75 73.04 76.97	67.44 54.13 68.97	68.97 69.67 72.12	68.75 68.04 71.97
80 - 60 - 20 - -20 - -20 - -40 - -60 - -80 -	-75 -50 -22	a) Fast	80 - 60 - 20 - 0 - -20 - -40 - -60 - -80 - 5	(b) F:	ast+KDCI	80 - 60 - 20 - 20 - -20 - -40 - -60 - -80 -	(c)	DeepInv 25 0 25 50 7	80 60 40 20 0 -20 -40 -60 -80	(d) Dec	epInv+KDCI	75

Figure 5. T-SNE results of vanilla and KDCI-based models performance on Tiny-ImageNet dataset. KDCI helps models obtain clearer clustering results, which show its strong positive impact.

		Ground Truth	Vanilla Fast	w/ KDCI
jeNet	Ż	albatross	missile	albatross
lmaç		manhole_cover	petri_dish	manhole_cover
Net		coral_reef	lawn_mower	coral_reef
y-Image	Ta	lakeside	alp	lakeside
Ë	1	seashore	lampshade	seashore

Figure 6. Qualitative results of the vanilla and KDCI-based version on ImageNet and Tiny-ImageNet.

nomenon *under-intervention*. The incomplete causal intervention would lead to gain reduction.

4.7. Qualitative Results

Further, we present qualitative results to further demonstrate the positive gains of KDCI over baseline methods.

Visualization Results. To intuitively show the help of KDCI to existing DFKD methods, we first visualize the student classification results with t-SNE [67] on the Tiny-ImageNet dataset. We reserve 100 classes of validating samples. From Figure 5, the KDCI-based versions (b)&(d) have fewer outliers and clearer clustering effects than the vanilla versions (a)&(c). These phenomena further confirm that our KDCI can well disentangle features from different

classes, thus improving existing methods' performance.

Case Study of Causal Intervention. We select representative examples from ImageNet and Tiny-ImageNet datasets to show differences in student predictions before and after the intervention. As shown in Figure 6, KDCI can eliminate the prediction offset caused by some misleading features to a certain extent. For example, students from the vanilla Fast misclassify "albatross" as "missile" or "coral_reef" as "lawn_mower" due to large patches of similar background colour, and misclassify "manhole_cover" as "petri_dish" or "seashore" as "lampshade" due to similar shape. Fortunately, KDCI can repair prediction shifts in the above cases.

5. Conclusion

This paper proposes a novel perspective from causal inference to handle the distribution shifts in the Data-Free Knowledge Distillation (DFKD) task. By customizing the causal graph according to the properties of the variables in the DFKD, we propose a Knowledge Distillation Causal Intervention (KDCI) framework to de-confound the adverse effect caused by the shifts between the substitution and original data. KDCI can be flexibly combined with most existing methods. Numerous experiments prove that KDCI can consistently help existing methods and provide an alternative causal intervention insight.

Acknowledgements

This work is supported by the Shanghai Engineering Research Center of AI & Robotics, Fudan University, China, the Engineering Research Center of AI & Robotics, Ministry of Education, China, and the Green Ecological Smart Technology School-Enterprise Joint Research Center.

References

- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceed*ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 770–778, 2016. 1, 5
- [2] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), pages 10012–10022, 2021.
- [3] Dingkang Yang, Yang Liu, Can Huang, Mingcheng Li, Xiao Zhao, Yuzheng Wang, Kun Yang, Yan Wang, Peng Zhai, and Lihua Zhang. Target and source modality co-reinforcement for emotion understanding from asynchronous multimodal sequences. *Knowledge-Based Systems*, page 110370, 2023.
- [4] Yang Liu, Dingkang Yang, Gaoyun Fang, Yuzheng Wang, Donglai Wei, Mengyang Zhao, Kai Cheng, Jing Liu, and Liang Song. Stochastic video normality network for abnormal event detection in surveillance videos. *Knowledge-Based Systems*, 280:110986, 2023.
- [5] Dingkang Yang, Kun Yang, Mingcheng Li, Shunli Wang, Shuaibing Wang, and Lihua Zhang. Robust emotion recognition in context debiasing. *arXiv preprint arXiv:2403.05963*, 2024.
- [6] Dingkang Yang, Shuai Huang, Haopeng Kuang, Yangtao Du, and Lihua Zhang. Disentangled representation learning for multimodal emotion recognition. In *Proceedings of the 30th* ACM International Conference on Multimedia (ACM MM), pages 1642–1651, 2022.
- [7] Dingkang Yang, Haopeng Kuang, Shuai Huang, and Lihua Zhang. Learning modality-specific and -agnostic representations for asynchronous multimodal language sequences. In *Proceedings of the 30th ACM International Conference on Multimedia (ACM MM)*, pages 1708–1717, 2022. 1
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 1
- [9] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4401–4410, 2019.
- [10] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. Advances in Neural Information Processing Systems (NeurIPS), 33:1877–1901, 2020.
- [11] Dingkang Yang, Kun Yang, Yuzheng Wang, Jing Liu, Zhi Xu, Rongbin Yin, Peng Zhai, and Lihua Zhang. How2comm: Communication-efficient and collaboration-pragmatic multiagent perception. Advances in Neural Information Processing Systems (NeurIPS), 36, 2024.
- [12] Zhaoyu Chen, Bo Li, Jianghe Xu, Shuang Wu, Shouhong Ding, and Wenqiang Zhang. Towards practical certifiable patch defense with vision transformer. In *Proceedings of*

the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 15148–15158, 2022.

- [13] Yang Liu, Jing Liu, Kun Yang, Bobo Ju, Siao Liu, Yuzheng Wang, Dingkang Yang, Peng Sun, and Liang Song. Ampnet: Appearance-motion prototype network assisted automatic video anomaly detection system. *IEEE Transactions* on Industrial Informatics, 2023. 1
- [14] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning* (*ICML*), pages 8748–8763. PMLR, 2021. 1
- [15] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- [16] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition (CVPR), pages 248–255. Ieee, 2009. 5, 15
- [17] Dingkang Yang, Shuai Huang, Zhi Xu, Zhenpeng Li, Shunli Wang, Mingcheng Li, Yuzheng Wang, Yang Liu, Kun Yang, Zhaoyu Chen, et al. Aide: A vision-driven multi-view, multimodal, multi-tasking dataset for assistive driving perception. In *Proceedings of the IEEE/CVF International Conference* on Computer Vision, pages 20459–20470, 2023.
- [18] Yuzheng Wang, Zhaoyu Chen, Dingkang Yang, Yang Liu, Siao Liu, Wenqiang Zhang, and Lizhe Qi. Adversarial contrastive distillation with adaptive denoising. In *ICASSP* 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 1–5. IEEE, 2023.
- [19] Zuhao Ge, Lizhe Qi, Yuzheng Wang, and Yunquan Sun. Zoom-and-reasoning: Joint foreground zoom and visualsemantic reasoning detection network for aerial images. *IEEE Signal Processing Letters*, 29:2572–2576, 2022.
- [20] Siao Liu, Zhaoyu Chen, Yang Liu, Yuzheng Wang, Dingkang Yang, Zhile Zhao, Ziqing Zhou, Xie Yi, Wei Li, Wenqiang Zhang, et al. Improving generalization in visual reinforcement learning via conflict-aware gradient agreement augmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 23436–23446, 2023. 1
- [21] Raphael Gontijo Lopes, Stefano Fenu, and Thad Starner. Data-free knowledge distillation for deep neural networks. *arXiv preprint arXiv:1710.07535*, 2017. 1, 2
- [22] Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al. Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531, 2(7), 2015. 1
- [23] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. arXiv preprint arXiv:1412.6550, 2014.
- [24] Yuzheng Wang, Zhaoyu Chen, Dingkang Yang, Pinxue Guo, Kaixun Jiang, Wenqiang Zhang, and Lizhe Qi. Out of thin

air: Exploring data-free adversarial robustness distillation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 5776–5784, 2024. 1

- [25] Paul R Burton, Madeleine J Murtagh, Andy Boyd, James B Williams, Edward S Dove, Susan E Wallace, Anne-Marie Tasse, Julian Little, Rex L Chisholm, Amadou Gaye, et al. Data safe havens in health research and healthcare. *Bioinformatics*, 31(20):3241–3248, 2015. 1
- [26] Hongxu Yin, Pavlo Molchanov, Jose M Alvarez, Zhizhong Li, Arun Mallya, Derek Hoiem, Niraj K Jha, and Jan Kautz. Dreaming to distill: Data-free knowledge transfer via deepinversion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8715–8724, 2020. 1, 2, 3, 15
- [27] Gongfan Fang, Jie Song, Xinchao Wang, Chengchao Shen, Xingen Wang, and Mingli Song. Contrastive model inversion for data-free knowledge distillation. arXiv preprint arXiv:2105.08584, 2021. 1, 2, 3, 5, 15
- [28] Gongfan Fang, Yifan Bao, Jie Song, Xinchao Wang, Donglin Xie, Chengchao Shen, and Mingli Song. Mosaicking to distill: Knowledge distillation from out-of-domain data. Advances in Neural Information Processing Systems (NeurIPS), 34:11920–11932, 2021. 1, 2, 3, 5, 15
- [29] Gongfan Fang, Kanya Mo, Xinchao Wang, Jie Song, Shitao Bei, Haofei Zhang, and Mingli Song. Up to 100x faster datafree knowledge distillation. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 36, pages 6597–6604, 2022. 1, 2, 3, 5, 12, 15
- [30] Kien Do, Thai Hung Le, Dung Nguyen, Dang Nguyen, Haripriya Harikumar, Truyen Tran, Santu Rana, and Svetha Venkatesh. Momentum adversarial distillation: Handling large distribution shifts in data-free knowledge distillation. Advances in Neural Information Processing Systems (NeurIPS), 35:10055–10067, 2022. 2, 3, 14
- [31] Hanting Chen, Yunhe Wang, Chang Xu, Zhaohui Yang, Chuanjian Liu, Boxin Shi, Chunjing Xu, Chao Xu, and Qi Tian. Data-free learning of student networks. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), pages 3514–3522, 2019. 2, 3, 5, 15
- [32] Hanting Chen, Tianyu Guo, Chang Xu, Wenshuo Li, Chunjing Xu, Chao Xu, and Yunhe Wang. Learning student networks in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6428–6437, 2021. 2, 3, 5, 12, 15
- [33] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 2, 5
- [34] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. Advances in Neural Information Processing Systems (NeurIPS), 30, 2017. 2
- [35] Judea Pearl. Causal inference in statistics: An overview. Statistics Surveys, 3:96–146, 2009. 2, 3, 4
- [36] Linda J Van Hamme and Edward A Wasserman. Cue competition in causality judgments: The role of nonpresentation of compound stimulus elements. *Learning and Motivation*, 25(2):127–151, 1994. 2

- [37] Madelyn Glymour, Judea Pearl, and Nicholas P Jewell. Causal inference in statistics: A primer. John Wiley & Sons, 2016. 2, 3
- [38] Paul Micaelli and Amos J Storkey. Zero-shot knowledge transfer via adversarial belief matching. *Advances in Neural Information Processing Systems (NeurIPS)*, 32, 2019. 2
- [39] Kartikeya Bhardwaj, Naveen Suda, and Radu Marculescu. Dream distillation: A data-independent model compression framework. arXiv preprint arXiv:1905.07072, 2019. 2, 5, 12
- [40] Liangchen Luo, Mark Sandler, Zi Lin, Andrey Zhmoginov, and Andrew Howard. Large-scale generative data-free distillation. arXiv preprint arXiv:2012.05578, 2020.
- [41] Yoojin Choi, Jihwan Choi, Mostafa El-Khamy, and Jungwon Lee. Data-free network quantization with adversarial knowledge distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 710–711, 2020.
- [42] Yuzheng Wang, Zuhao Ge, Zhaoyu Chen, Xian Liu, Chuangjia Ma, Yunquan Sun, and Lizhe Qi. Explicit and implicit knowledge distillation via unlabeled data. In *ICASSP* 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 1–5. IEEE, 2023. 2
- [43] Yuzheng Wang, Zhaoyu Chen, Jie Zhang, Dingkang Yang, Zuhao Ge, Yang Liu, Siao Liu, Yunquan Sun, Wenqiang Zhang, and Lizhe Qi. Sampling to distill: Knowledge transfer from open-world data. arXiv preprint arXiv:2307.16601, 2023. 2
- [44] Hal R Varian. Causal inference in economics and marketing. Proceedings of the National Academy of Sciences, 113(27):7310–7315, 2016. 2
- [45] E Michael Foster. Causal inference and developmental psychology. *Developmental Psychology*, 46(6):1454, 2010. 2
- [46] Tan Wang, Jianqiang Huang, Hanwang Zhang, and Qianru Sun. Visual commonsense r-cnn. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 10760–10770, 2020. 3
- [47] Yingjie Chen, Diqi Chen, Tao Wang, Yizhou Wang, and Yun Liang. Causal intervention for subject-deconfounded facial action unit recognition. arXiv preprint arXiv:2204.07935, 2022.
- [48] Xu Yang, Hanwang Zhang, Guojun Qi, and Jianfei Cai. Causal attention for vision-language tasks. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 9847–9857, 2021.
- [49] Yang Liu, Zhaoyang Xia, Mengyang Zhao, Donglai Wei, Yuzheng Wang, Siao Liu, Bobo Ju, Gaoyun Fang, Jing Liu, and Liang Song. Learning causality-inspired representation consistency for video anomaly detection. In *Proceedings* of the 31st ACM International Conference on Multimedia (ACM MM), pages 203–212, 2023.
- [50] Dingkang Yang, Dongling Xiao, Ke Li, Yuzheng Wang, Zhaoyu Chen, Jinjie Wei, and Lihua Zhang. Towards multimodal human intention understanding debiasing via subjectdeconfounding. arXiv preprint arXiv:2403.05025, 2024.
- [51] Xiang Deng and Zhongfei Zhang. Comprehensive knowledge distillation with causal intervention. *Advances in Neu-*

ral Information Processing Systems (NeurIPS), 34:22158–22170, 2021. 3

- [52] Kaihua Tang, Yulei Niu, Jianqiang Huang, Jiaxin Shi, and Hanwang Zhang. Unbiased scene graph generation from biased training. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition (CVPR), pages 3716–3725, 2020. 3
- [53] Teng Sun, Wenjie Wang, Liqiang Jing, Yiran Cui, Xuemeng Song, and Liqiang Nie. Counterfactual reasoning for out-ofdistribution multimodal sentiment analysis. arXiv preprint arXiv:2207.11652, 2022.
- [54] Chen Qian, Fuli Feng, Lijie Wen, Chunping Ma, and Pengjun Xie. Counterfactual inference for text classification debiasing. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 5434–5445, 2021.
- [55] Dingkang Yang, Mingcheng Li, Dongling Xiao, Yang Liu, Kun Yang, Zhaoyu Chen, Yuzheng Wang, Peng Zhai, Ke Li, and Lihua Zhang. Towards multimodal sentiment analysis debiasing via bias purification. arXiv preprint arXiv:2403.05023, 2024. 3
- [56] Judea Pearl. *Causality*. Cambridge University Press, 2009.3
- [57] Yulei Niu, Kaihua Tang, Hanwang Zhang, Zhiwu Lu, Xian-Sheng Hua, and Ji-Rong Wen. Counterfactual vqa: A causeeffect look at language bias. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (CVPR), pages 12700–12710, 2021. 3
- [58] Dingkang Yang, Zhaoyu Chen, Yuzheng Wang, Shunli Wang, Mingcheng Li, Siao Liu, Xiao Zhao, Shuai Huang, Zhiyan Dong, Peng Zhai, and Lihua Zhang. Context deconfounded emotion recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19005–19015, June 2023. 3
- [59] Judea Pearl et al. Models, reasoning and inference. Cambridge, UK: CambridgeUniversityPress, 19:2, 2000. 3
- [60] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International Conference on Machine Learning (ICML)*, pages 2048–2057. PMLR, 2015.
- [61] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015. 5
- [62] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556, 2014. 5
- [63] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. arXiv preprint arXiv:1605.07146, 2016. 5
- [64] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. arXiv preprint arXiv:1710.09412, 2017. 5
- [65] Borui Zhao, Quan Cui, Renjie Song, Yiyu Qiu, and Jiajun Liang. Decoupled knowledge distillation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 11953–11962, 2022. 5, 12

- [66] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning (ICML)*, pages 448–456. pmlr, 2015. 6
- [67] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(11), 2008. 8
- [68] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. Advances in Neural Information Processing Systems (NeurIPS), 32, 2019. 12
- [69] Kuluhan Binici, Nam Trung Pham, Tulika Mitra, and Karianto Leman. Preventing catastrophic forgetting and distribution mismatch in knowledge distillation via synthetic data. In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), pages 663–671, 2022. 14
- [70] Kuluhan Binici, Shivam Aggarwal, Nam Trung Pham, Karianto Leman, and Tulika Mitra. Robust and resourceefficient data-free knowledge distillation by generative pseudo replay. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 36, pages 6089–6096, 2022. 14
- [71] Gaurav Patel, Konda Reddy Mopuri, and Qiang Qiu. Learning to retain while acquiring: Combating distribution-shift in adversarial data-free knowledge distillation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 7786–7794, 2023. 14

De-confounded Data-free Knowledge Distillation for Handling Distribution Shifts

Supplementary Material

In this supplementary material, we provide more details of our method, organized as follows:

- In Section A, we provide the detailed training settings and illustrate how KDCI combines with existing DFKD methods, and show the algorithm process, corresponding to Section 4.3 of the main body.
- In Section B, we qualitatively assess students' learning progress about vanilla DFKD methods and their KDCIbased version to verify the positive effect of KDCI on the existing DFKD method.
- In Section C, we analyze the possible reasons for the difference in performance improvement, corresponding to Section 4 of the main body.
- In Section D, we provide more observable visualization results as more sufficient evidence, corresponding to Section 4.7 of the main body.
- In Section E, we discuss the significant differences between our KDCI and other methods focusing on data distribution.
- In Section F, we discuss the broader impact and potential limitations.
- In Section G, we provide the detailed experimental settings for the used baseline methods, corresponding to Section 4.2 of the main body.

A. Additional Training Details & Algorithm Process of Combining KDCI with Existing DFKD Methods

A.1. Training Details

We provide the detailed experimental settings for our KDCI framework. Our KDCI and reproducible methods are implemented through PyTorch [68]. All models are trained on RTX 3090 GPUs. For CIFAR-10 and CIFAR-100, all training settings (e.g., loss function, optimizer, batch size, learning rate, etc) of the reported methods are consistent with the released codebase. The results are shown in Table 1 of the main body. For Tiny-ImageNet, initially, we try to find a unified teacher model for the Tiny-ImageNet dataset in open-sourced projects. However, one problem is that the teacher model pre-trained on Tiny-ImageNet seems confidential, so finding an open-source unified model is difficult. In this case, we train the unified renset-34 teacher model for 200 epochs on the original training data. During the teacher's training, we use the SGD optimizer with the momentum as 0.9, weight decay as 5e-4, the batch size as 128, and cosine annealing learning rate with an initial value of 0.1. The teacher model can converge without addi-

Algorithm 1 Training process of generation-based methods combined with our KDCI

- **Input:** A pre-trained teacher model T, a generator g, a student model S, distillation epochs T, batch size N_m , the iterations of generator g in each epoch Tg, the iterations of student f_s in each epoch Ts, the confounder size N.
- 1: **for** epoch = [1, ..., T] **do**
- 2: // Generation stage
- 3: **for** generator iterations $= [1, \ldots, Tg]$ **do**
- 4: Randomly sample noises and labels (z, y)
- 5: Synthesize a mini-batch training data X = g(z, y)
 - Update generator g with the generator loss
- 7: end for

6:

- 8: Synthesize training data X = g(z, y). Obtain the prediction feature $M = \{m_j \in \mathbb{R}^d\}_{i=1}^{N_m}$
- 9: Prototype clustering for M. Calculate the number of the prediction features in *i*-th cluster N_i , the feature cluster $\sum_{k=1}^{N_i} m_k^i$ and the subcenter
- $\boldsymbol{z}_{i} = \frac{1}{N_{i}} \sum_{k=1}^{N_{i}} m_{k}^{i}.$ 10: Construct a confounder dictionary $\boldsymbol{Z} = [\boldsymbol{z}_{1}, \boldsymbol{z}_{2}, \dots, \boldsymbol{z}_{N}]$ and calculate the prototype proportion $P_{s}(\boldsymbol{z}_{i}) = N_{i}/N_{m}$
- 11: // Distillation stage
- 12: **for** student iterations $= [1, \ldots, T_s]$ **do**
- 13: Synthesize training data X = g(z, y). Get models's predictions T(X) and S(X)
- 14: Calculate the prior information: $F(\boldsymbol{z}) = \sum_{i=1}^{N} \lambda_i \boldsymbol{z}_i P_s(\boldsymbol{z}_i)$
- 15: Compensate the student's predictions: $S'(X) = \phi(S(X), F(z))$
- 16: Update the student S with $KD\langle T(X), S'(X) \rangle$
- 17: end for
- 18: end for
- **Output:** The student model *S*.

tional tuning. Based on this pre-trained teacher, we train all students for 200 epochs. For the student, we use the SGD optimizer with the momentum as 0.9, the weight decay as 1e-4, the batch size as 256, the cosine annealing learning rate with an initial value of 0.2 for Fast [29], and 0.1 for DeepInv [39] & DFND [32]. The results are shown in Table 2 of the main body. For ImageNet, We choose the same pre-trained resnet-50 model with [65] and unify the teacher model of different baseline methods. For Fast, we test directly on the open-source project. For DeepInv, we reproduce the corresponding results with the specified backbone pair. For DFND, we select 600k samples from the unlabeled FlickerlM dataset. The teacher's backbone is different from the original paper. The different backbones may cause the results we reproduce to differ from the original paper. The results are shown in Table 1 of the supplementary material.



Figure 7. The test accuracy on Tiny-ImageNet dataset across different local training epochs $E = \{10, 20, \dots, 200\}$. Our KDCI framework improves the performance of baselines consistently.

For the implementation of our KDCI, the hidden dimension d_n is set to 256. And d_h equals the hidden dimension d and the number of classes. By default, $\phi(\cdot)$ uses feature addition. For various baseline methods, the settings are shown in Section G of the supplementary material.

A.2. Algorithm Process

In the existing DFKD task, the generation-based and sampling-based method processes are different. Therefore, the way KDCI combines these methods and the hyperparameter settings are also slightly different. For the generation-based process, the generator and student models are updated alternately, which means the student's training data is updated in each epoch. We use a mini-batch of synthetic training data to construct the confounder dictionary, and the dictionary will be updated as the generator is updated. For the sampling-based process, existing methods select unlabeled data according to the preferences of the teacher model. Then, the student relies on these unlabeled data for data-based knowledge distillation training. We use all sampled data to construct the confounder dictionary. During subsequent student training, the dictionary is fixed. For a clearer understanding, we describe the above process as Algorithm 1 and 2, respectively.

B. Vanilla DFKD Methods vs. Their KDCI**based Versions**

In the main body, we have compared the quantitative results of vanilla DFKD methods and their KDCI-based versions. To observe the positive effect of KDCI on the existing DFKD methods more clearly, we visualize the student's test accuracy on the Tiny-ImageNet dataset. The results are shown in Figure 7. KDCI can consistently help students from the beginning of training to the end, which verifies its effectiveness.

Algorithm 2 Training process of sampling-based methods combined with our KDCI

- Input: A pre-trained teacher model T, a student model S, unlabeled training dataset $D = \{x_j\}_{j=1}^n$, distillation epochs T, batch size m, number of batches M, the number of sampled data N_m , the confounder size N.
- 1: // Sampling stage
- 2: Sample the training data $\{x_j\}_{j=1}^{N_m}$ from D. Obtain the prediction feature set $M = \{m_j \in \mathbb{R}^d\}_{i=1}^{N_m}$
- 3: Prototype clustering for M. Calculate the number of the prediction features in i-th cluster N_i , the feature cluster $\sum_{k=1}^{N_i} m_k^i$ and the subcenter $\boldsymbol{z}_i = \frac{1}{N_i} \sum_{k=1}^{N_i} m_k^i$.
- 4: Construct a confounder dictionary $\mathbf{Z} = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N]$ and calculate the prototype proportion $P_s(z_i) = N_i/N_m$
- 5: // Distillation stage
- for epoch = $[1, \ldots, T]$ do 6:
- for mini-batch = $[1, \ldots, M]$ do 7:
- 8: Sample a mini-batch training data: $\mathbf{X} = \{x_i\}_{i=1}^m$ from $\{x_j\}_{j=1}^{N_m^m}$ Get teacher and student predictions $\mathbf{T}(\mathbf{X})$ and $\mathbf{S}(\mathbf{X})$
- 9:
- 10: Calculate the prior information:
- $F(\boldsymbol{z}) = \sum_{i=1}^{N} \lambda_i \boldsymbol{z}_i P_s(\boldsymbol{z}_i)$
- Compensate the student's predictions: 11:
- $S'(X) = \phi(S(X), F(z))$

12: Update the student **S** with
$$KD\langle T(\mathbf{X}), S'(\mathbf{X})\rangle$$

- end for 13:
- 14: end for
- Output: The student model S.

C. Analyses of Difference in Performance Improvements

Judging from the experimental results, KDCI has different gains for different DFKD methods on different datasets. We think such observations arise from various factors.

• By default, we choose the teacher model itself to extract the confounding dictionary. The prediction feature set provided by teachers of different backbones has different expressiveness, which affects the compensation degree of backdoor adjustment for bias during the causal intervention. The tests in Lines 513-531 and Table. 5 of the main body also verify this conclusion.

- The degree of distribution shift of synthetic data on distinct datasets is different. More complex datasets may degrade the generation quality for generation-based methods, resulting in more significant distribution shifts. KDCI tends to be more effective for more sophisticated datasets.
- Different baseline methods with different training losses are influential. Observations such as Section 4.4 of the main body suggest that methods that already incorporate prior likelihood knowledge of the data may weaken the KDCI gain.
- In addition, there may be many underlying factors. Nevertheless, KDCI, as a model-agnostic general framework, has promising and competitive improvements and gains for various models as a whole. We believe that a deeper exploration of the relevant mechanisms is a promising perspective. For this topic, we leave it to future work.



Figure 8. Qualitative results of the vanilla and KDCI-based version on CIFAR-10, CIFAR-100, ImageNet, and Tiny-ImageNet.

D. More Visual Evidence

To further verify the effectiveness, we provide more case studies of causal intervention. As shown in Figure 8, we visualize some test instances corrected by our KDCI compared to the vanilla version (Fast) on four kinds of datasets (*i.e.*, CIFAR-10, CIFAR-100, ImageNet, and Tiny-ImageNet). The vanilla version sometimes confuses some test instances due to shape or color. Our KDCI can repair these prediction shifts to enhance student performance.

E. Discussion with Other Works that Address Distribution Shifts

Several DFKD works already address distribution shifts in adversarial contexts [30, 69–71]. The works reveal distribution shift issues in the DFKD task from different aspects, but our method is significantly different from these works. Specifically, the differences between our KDCI and others are as follows:

- **Applicability.** These existing works tacitly use the same motivation, *i.e.*, as the generator gets updated, the distribution of synthetic data will change, causing the student to forget the knowledge it acquired at previous steps. However, such motivation does not apply to sampling-based methods. After selecting the training samples, they will not change during the entire student training process. Our motivation comes from the observed distribution shifts between the substitution data and can cover the two methods mentioned.
- Economy. Existing methods often rely on substantial additional computational and storage costs, *e.g.*, the need to store and maintain an additional dynamic collection of generated samples [69], the need for additional generator architectures to memorize knowledge of past generated data (an additional Variational Autoencoder (VAE) [70] or Exponential Moving Average generator [30]), and additional memory bank or additional loss calculation and gradient update [71]. In contrast, our method only needs to compute and store a small number of matrix computation results. Compared with the update of the models, the computational cost of the clustering process is basically negligible.
- **Plug-and-play.** Existing works are to propose new methods. Undoubtedly, these methods can provide a potential reference for other DFKD methods, but whether they can be easily combined with existing DFKD methods and improve overall performance is still unknown. Our proposed technique is model-agnostic, as a plug-and-play paradigm that integrates well with existing works. A large number of experiments have proved this conclusion.

F. Further Discussion

F.1. Broader Impact

The positive impact of this work: the proposed KDCI module can suppress the distribution shifts between the substitution and original data in the DFKD task, preventing the potential discrimination of the student's learning. While the pre-trained model for extracting prior knowledge uses the teacher itself, our method does not require additional dependencies and auxiliary information. The negative impacts of this work: students may be forced to identify minority groups for malicious purposes with customized biased teacher models. Therefore, we have to make sure that the DFKD technique is used for the right purpose.

F.2. Limitations

Since there are countless methods with insights for the DFKD task, other ways of classifying forms may also be reasonable. In this paper, we simply divide the source of the substitution data into generation-based and sampling-based methods. Similarly, it is impossible to cover all DFKD methods, so only open-source and representative methods are selected as the baseline. Nevertheless, the existing performance improvement is enough to prove the positive impact of KDCI on students.

In addition, since what we propose is a framework rather than a specific method, the test on the effectiveness of KDCI relies on the experimental setting of the existing DFKD methods. Currently, the mainstream open-source DFKD methods rarely use real-life medical or facial datasets for testing, so we only follow the mainstream experimental settings. Following the consensus of peers is necessary to increase the impact of our work. In this work, we select datasets that are widely used and accepted by the vast majority of DFKD methods. Following previous data paradigms is beneficial for acceptance by the relevant research community and enhances the persuasiveness of our method.

G. Experimental Setup of the Baseline DFKD Methods

DAFL. DAFL [31] is a data-free generation method. We keep the generator loss from the original as: $\mathcal{L}_{GEN} = \mathcal{L}_{oh} + \alpha \mathcal{L}_a + \beta \mathcal{L}_{ie}$. The knowledge distillation loss is: $\mathcal{L}_{KD} = D_{KL}(\mathcal{N}_S(x), \mathcal{N}_T(x))$. Following the original settings, we set $\alpha = 1e - 3$, $\beta = 20$. We use SGD with the weight decay of 5e - 4, the momentum of 0.9, and the initial learning rate set as 0.1.

Fast. Fast [29] is a fast data-free generation method via feature sharing. We keep the generator loss from the original as: $\mathcal{L}_{GEN} = \alpha \mathcal{L}_{cls} + \beta \mathcal{L}_{adv} + \gamma \mathcal{L}_{feat}$. The knowledge distillation loss is: $\mathcal{L}_{KD} = D_{KL}(\mathcal{N}_S(x), \mathcal{N}_T(x))$. We set $\alpha = 0.4, \beta = 1.1$, and $\gamma = 10$, which are the same as

the original settings. We use the Adam Optimizer with a learning rate of 1e - 3 to update the generator and the SGD optimizer with a momentum of 0.9 and a learning rate of 0.1 for student training.

CMI. CMI [27] is a model inversion method with contrastive learning. We keep the generator loss from the original as: $\mathcal{L}_{GEN} = \alpha \mathcal{L}_{bn} + \beta \mathcal{L}_{cls} + \gamma \mathcal{L}_{adv} + \delta \mathcal{L}_{cr}$. The knowledge distillation loss is: $\mathcal{L}_{KD} = D_{KL}(\mathcal{N}_S(x), \mathcal{N}_T(x))$. We set $\alpha = 1$, $\beta = 0.5$, $\gamma = 0.5$, and $\delta = 0.8$. We use the Adam Optimizer with a learning rate of 1e - 3 to update the generator and the SGD optimizer with a momentum of 0.9 and a learning rate of 0.1 for student training.

DeepInv. DeepInv [26] is a model inversion method that combines prior knowledge and adversarial training. We keep the inversion loss from the original as: $\mathcal{L}_{GEN} = \alpha_{tv}\mathcal{R}_{tv} + \alpha_{l2}\mathcal{R}_{l2} + \alpha_{f}\mathcal{R}_{feature} + \alpha_{c}\mathcal{R}_{compete}$. The knowledge distillation loss is: $\mathcal{L}_{KD} = D_{KL}(\mathcal{N}_{S}(x), \mathcal{N}_{T}(x))$. We set $\alpha_{tv} = 2.5e - 5$, $\alpha_{l2} = 3e - 8$, $\alpha_{f} = 0.1$ and $\alpha_{c} = 10$, which are the same as the original setting. Besides, we set the number of iterations as 1000 and use Adam for optimization with a learning rate of 0.05.

DFND. DFND [32] is a sampling-based method using open-world unlabeled data as the substitution data. Following the original, we select 600k data with the highest teacher confidence from the ImageNet dataset [16] as the sampled data and resize them to the resolution of the corresponding dataset. We use the same noisy distillation loss $\mathcal{L}_{KD} = \mathcal{H}_{CE}(Q(\mathcal{N}_S(x)), \hat{y}) + \lambda D_{KL}(\mathcal{N}_S(x), \mathcal{N}_T(x))$, and λ is set as 4. The student network is optimized using SGD and the initial learning rate is set as 0.1 Weight decay and momentum are set as 5e - 4 and 0.9, respectively.

Mosaick. Mosaick [28] is a sampling-based method using out-of-domain (OOD) unlabeled data as the substitution data. We select 600k data with the lowest teacher confidence from the ImageNet dataset [16] as the OOD data. Following the original settings, we use Adam for optimization, with hyper-parameters lr = 1e - 3, $\beta_1 = 0.5$, and $\beta_2 = 0.999$ for the generator and discriminator. The distillation loss is $\mathcal{L}_{KD} = \lambda D_{KL} - \lambda \mathcal{R}(G, D, T)$ The student network is optimized using SGD, and the initial learning rate is set as 0.1. Weight decay and momentum are set as 1e - 4 and 0.9, respectively.