# Data-free Knowledge Distillation for Fine-grained Visual Categorization

Renrong Shao[1], Wei Zhang[1][*], Jianhua Yin[2], Jun Wang[1][*]

[1]School of Computer Science and Technology, East China Normal University
[2]School of Computer Science and Technology, Shandong University

{roryshaw6613,zhangwei.thu,jhyinmail, wongjun}@gmail.com

## Abstract

*Data-free knowledge distillation (DFKD) is a promising approach for addressing issues related to model compression, security privacy, and transmission restrictions. Although the existing methods exploiting DFKD have achieved inspiring achievements in coarse-grained classification, in practical applications involving fine-grained classification tasks that require more detailed distinctions between similar categories, sub-optimal results are obtained. To address this issue, we propose an approach called DFKD-FGVC that extends DFKD to fine-grained visual categorization (FGVC) tasks. Our approach utilizes an adversarial distillation framework with attention generator, mixed high-order attention distillation, and semantic feature contrast learning. Specifically, we introduce a spatial-wise attention mechanism to the generator to synthesize fine-grained images with more details of discriminative parts. We also utilize the mixed high-order attention mechanism to capture complex interactions among parts and the subtle differences among discriminative features of the fine-grained categories, paying attention to both local features and semantic context relationships. Moreover, we leverage the teacher and student models of the distillation framework to contrast high-level semantic feature maps in the hyperspace, comparing variances of different categories. We evaluate our approach on three widely-used FGVC benchmarks (Aircraft, Cars196, and CUB200) and demonstrate its superior performance. Code is available at https://github.com/RoryShao/DFKD-FGVC.git*

## 1. Introduction

Fine-grained visual categorization (FGVC) aims at distinguishing subcategories from father categories, e.g., subcategories of birds [43], aircraft [29], and cars [23]. It has long been considered a more challenging issue than traditional image classification due to the subtle inter-class and large intra-class variations [42]. To distinguish subtle diversities, the current approaches commonly exploit deeper neural networks with elaborate designs [50, 26, 53] to excavate the discriminative features effectively. Inevitably, the network becomes more and more complex, which leads to another problem, i.e., complicated networks are not easily deployed on embedded or mobile devices. Besides, the training data of released pre-trained models are often unavailable due to transmission, privacy, or legal issues. For example, pre-trained models commonly need a large amount of data such as ImageNet [24]. If the data is transmitted directly, a large amount of bandwidth is consumed. Moreover, some sensitive data such as e-commerce items or medical data are usually not directly accessible to the public due to intellectual property rights or privacy protection considerations. To obtain a lightweight model, recent research has made significant progress, including pruning [25], quantization [49, 27], and knowledge distillation [16]. Among them, knowledge distillation (KD) is a popular and effective paradigm for model compression and knowledge transfer [16]. It works by transferring knowledge from a cumbersome teacher network to a lightweight student network. Thanks to this separable architecture, it can also be used to solve privacy protection in data-free scenarios, which is called data-free knowledge distillation (DFKD) [5] or zero-shot knowledge distillation (ZSKD) [33].

Fortunately, a series of DFKD methods have been proposed [5, 33, 30, 11, 47, 12]. The existing approaches can be divided into two paradigms. The first paradigm is based on the category distribution, which exploits the out distribution of teacher and student to optimize the student and generator, e.g., DFAL [5], ZSKT [30], DFAD [11], ZSKD [33]. Such a paradigm commonly fails to generate realistic samples due to the lack of semantic-related information, especially when it comes to complex samples. The second paradigm is based on prior distribution, which exploits the prior information (i.e., BatchNorm) to optimize synthetic images for distillation, e.g., MAD [10], CMI [12], DFQ [8],

ADI [47]. This paradigm can produce realistic features and, therefore, gives the student a noticeable improvement.

Although the existing methods have achieved inspiring achievements in coarse-grained classification, in practical applications, sub-optimal results are achieved due to the subtle variations widely found in different scenarios. The main reasons for this situation are as follows: Firstly, for FGVC tasks, the variances of the same category are more prominent than that of coarse-grained classification due to different factors, such as viewing angles, lighting, backgrounds, occlusion, etc. Secondly, compared to coarse-grained categories, the feature discrepancies of different categories in FGVC are not obvious. Besides, in the data-free scenario, the model can not access the raw data directly. For synthesized images, it is difficult for the teacher model to capture the subtle variances of discriminative features. To our best knowledge, there are still no specialized data-free studies on fine-grained DFKD. Therefore, this inspires us to explore this issue and tackle this task in a data-free scenario.

In this paper, we tackle this issue by extending DFKD to fine-grained visual classification (FGVC) tasks and propose an approach named DFKD-FGVC, which is achieved by exploiting the adversarial distillation framework with attention generator, mixed high-order attention distillation (MHAD) and semantic feature contrast learning (SFCL). Concretely, as shown in Fig. 1, to promote the generator to synthesize more fine-grained images, we exploit the generator with spatial-wise attention, which can help the generator synthesize the images with more details of discriminative parts. Then, to fully mine the knowledge of discriminative features for student, we exploit the mixed high-order attention mechanism to capture complex interactions among parts and the subtle differences among discriminative features of the fine-grained categories, paying attention to both local features and semantic context relationships. Besides, to compare variances of different categories, we skillfully exploit the teacher and student model of distillation framework to contrast semantic feature maps in the hyperspace. To verify our approach, massive experiments are conducted on three fine-grained benchmarks, such as Aircraft, Cars196, and CUB200 to evaluate the effectiveness of our approach.

In a nutshell, our contributions are four-fold: 1) We are the first to propose an approach for FGVC in the data-free distillation scenario, which aims to optimize the entire generation and distillation stages to focus on discriminative features. 2) To synthesize more fine-grained images for adversarial distillation, we employ the generator with spatial-wise attention, which motivates the generator to synthesize the images with more details of discriminative features. 3) Particularly, to effectively mine the potential semantic features and contextual relationships of the fine-grained categories, we provide two strategies, namely, MHAD and SFCL, both of which can promote the performance of DFKD from different dimensions. 4) Extensive experiments demonstrate the effectiveness of our approach in the data-free setting, which achieves state-of-the-art performance on the mainstream FGVC benchmark datasets.

## 2. Related Works

### 2.1. Fine-Grained Visual Categorization

Fine-grained visual classification (FGVC) [43, 29, 23] is much more challenging than traditional classification tasks due to the inherently subtle intra-class object variations [42, 18]. Benefiting from the recent development of neural networks, recent studies have moved from strongly supervised information with extra annotations such as bounding boxes [2, 48, 18] to weakly-supervised conditions with only category labels [51, 13, 41]. Current methods on FGVC can be roughly divided into localization-based methods [13, 41] and attention-based methods [3, 19, 31]. The core for solving FGVC is to learn the discriminative features of objects in images. However, current approaches tackle this problem in the data-driven setting, few approaches consider this problem in the data-free setting. Therefore, different from the above studies, we explore the FGVC tasks in the novel aspect of the data-free distillation scenario.

### 2.2. Attention Mechanism

The attention mechanism stems from human vision, which exploits a sequence of partial glimpses and selectively focuses on salient parts to capture visual structure better. In the field of computer vision, attention mechanism [40, 17, 44, 34] are mainly exploited to capture essential information in various tasks such as pedestrian re-identification [46, 20, 45], FGVC [3, 31], etc. For example, [40] proposes the residual attention network for large-scale classification tasks. Then Hu et al. [17] exploit a squeeze-and-excitation (SE) block to compute channel-wise attention. CBAM [44] infers attention maps along two separate dimensions, i.e., channel and spatial. Similar to [44], BAM [34] also exploits the 3D attention map inference into channel and spatial. In terms of tasks, spatial attention is well-suited to dense prediction tasks such as semantic segmentation and object detection [14], while channel attention is a good choice for image classification. However, only exploiting spatial attention or channel attention is coarse, we can not capture the high-order and complex interactions among parts [4]. Therefore, in our data-free framework, we empirically exploit the spatial attention for our generator and the mixed high-order attention for distillation.

### 2.3. Data-free Distillation

Data-free Distillation has become a hot topic in recent years, mainly due to privacy protection [28]. It ex-
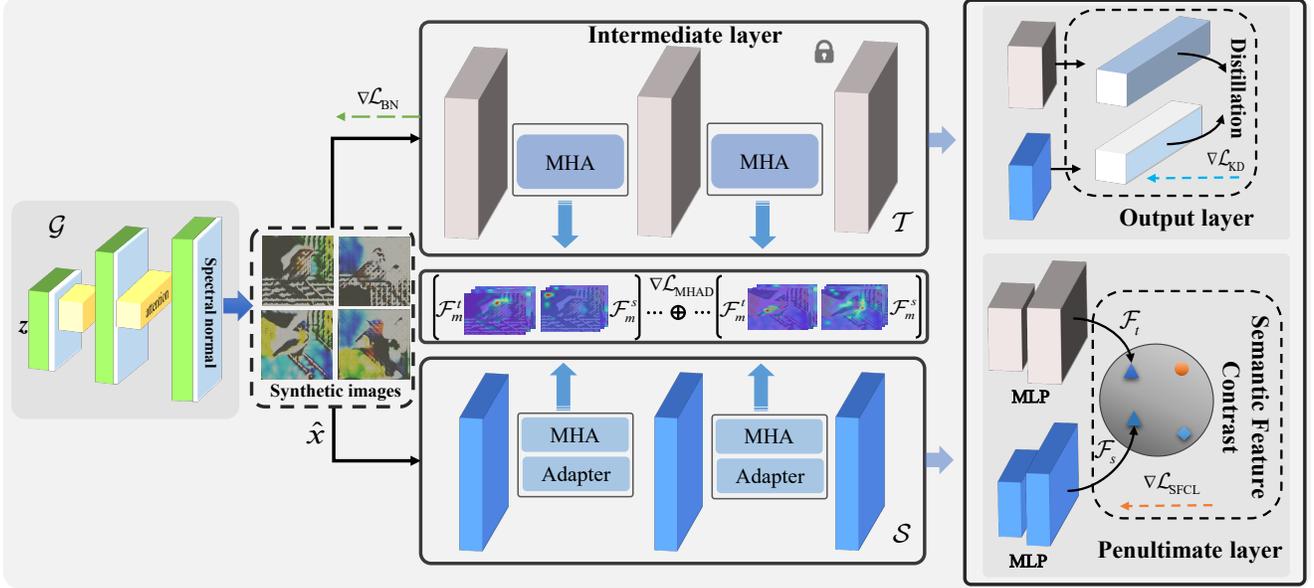
Figure 1. The whole framework of our approach. The **left**: The spatial attention module is plugged into each block of generator $\mathcal{G}$, which aims to focus on fine-grained semantic information from the whole process of noise $z$ to images $\hat{x}$. The **intermediate**: At each block of teacher and student, the feature maps are extracted by the mixed high-order attention module to achieve MHAD. The **right**: In the penultimate layer, exploiting the MLP to map the high-level semantic features of teacher and student to a common hyperspace and compare the variances by SFCL.

ploits synthesized alternative samples to solve the dilemma that model can not directly access the original data and makes gratifying achievements in the task of classification [47, 5, 11, 12]. For example, ADI [47] utilizes batch normalization statistics (BNS) of the pre-trained teacher to optimize the noise to synthesize high-fidelity images for KD. CMI [12] exploits the local and global contrast of samples to optimize the generator diversity. This kind of method ordinarily can synthesize more realistic images and achieve relatively better performance. DFAL [5] adopts a generator to synthesize images, and then the student learns the knowledge from the teacher by distillation. ZSKT [30] exploits the adversarial distillation to transfer the knowledge from teacher to student by KL and spatial attention, while DFAD [11] only utilizes the MAE loss to fit the output distribution of the teacher. All kinds of the above methods can achieve relatively inspiring achievements in coarse-grained classification, and there is no specific research on FGVC. Motivated by this, we conduct the study for data-free fine-grained distillation.

## 3. Preliminary

Our approach follows the basic thinking of DFKD, as depicted in Fig. 1. First, a generator $\mathcal{G}$ is employed to synthesize a batch of images from noise $z \sim \mathcal{N}(0,1)$, $\mathcal{G}(z) \to \hat{x}$, $\mathcal{B} = \{\hat{x}_1, \hat{x}_2, ..., \hat{x}_n\}, n \in \{1, ..., N\}$, where N is the batch size. Then the synthesized image $\hat{x}$ is input to the pre-trained teacher $\mathcal{T}$ and student $\mathcal{S}$ to support their distillation.

Finally, the generator $\mathcal{G}$ is optimized by adversarial distillation.

**Data-free Adversarial Distillation.** Essentially, Data-free adversarial distillation is a robust minimax optimization problem [1], which encourages the generator to minimize the possible loss for a worst-case scenario (maximum loss) through adversarial training under data uncertainty. In the data-free scenario, it can be denoted as

$$\min_{\mathcal{S}} \max_{\mathcal{G}} \left\{ \mathbb{E}_{p(z)} \left[ \mathcal{D}(\mathcal{T}(\mathcal{G}(z)), \mathcal{S}(\mathcal{G}(z))) \right] - \delta \mathcal{L}_{\mathcal{G}} \right\}, \quad (1)$$

where $\mathcal{D}$ represents the discrepancy measure, which normally exploits the Kullback-Leibler (KL) divergence as an optimization term. $\delta \geq 0$ is the balance factor, and $\mathcal{L}_{\mathcal{G}}$ is the optimization term of generator $\mathcal{G}$.

**Knowledge Distillation.** According to the principle of classic knowledge distillation [16], the soft output of the network (a.k.a. probability distribution) implies the similarity between the current sample and other categories. Therefore, traditional methods [8, 10, 33] usually adopt the KL Divergence to measure the difference between the two distributions of teacher and student. The probability distribution distillation can be formulated as

$$\mathcal{L}_{\text{KD}} = \mathbb{E}_{\hat{x}} \left[ D_{\text{KL}} \left( \sigma(\mathcal{S}(\hat{x})) \| \sigma(\mathcal{T}(\hat{x})) \right) \right], \quad (2)$$

where $D_{\text{KL}}$ represents the Kullback-Leibler (KL) divergence, and $\sigma$ is the softmax operation.

**Prior Information Regularization.** Prior information regularization aims to regularize the feature distribution of syn-
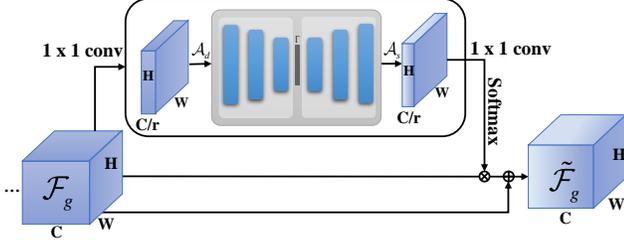
Figure 2. The spatial attention module of the generator, in which $\otimes$ denotes the element-wise multiplication and $\oplus$ denotes the element-wise addition.

thesized images by prior distribution information, i.e., mean $\mu$ and variance $\sigma^2$ of BatchNorm [47], which motivates the synthetic samples to approach the distribution of the original samples.

$$\mathcal{L}_{\mathrm{BN}} = \min_{\mathcal{G}} \sum_l \|\mu_l - \mu_l(\mathcal{G}(z))\|_2 + \left\|\sigma_l^2 - \sigma_l^2(\mathcal{G}(z))\right\|_2 , \quad (3)$$

where $l$ donates the $l^{th}$ BatchNorm layer of the teacher model, $\mu$ and $\sigma^2$ are the batch-wise mean and variance, respectively.

## 4. Proposed Approach

### 4.1. Discriminative Feature Synthesis

In the DFKD framework, it is common to exploit a generator to assist in generating alternative samples for coarse-grained classification. However, directly applying it to synthesize fine-grained samples often does not yield desirable discriminative features. This is because the conventional generator cannot focus on subtle discriminative features, which decreases the ability of teacher to extract representation from various semantic parts and thus hampers the effectiveness of the distillation. Differing from the traditional approaches, in our framework, we employ a DCGAN [35, 38] generator with the attention module to increase the representation ability of features and tell the generator where to focus. Inspired by preceding attention works such as CBAM [44] and CBM [34], which stacks channel attention and spatial attention in series, we exploit the attention mechanism in our approach. However, unlike the prior approaches, we implement the attention by the encoder-decoder manner, thinking that the non-linear convolution can pay attention to context knowledge of features, which is more suitable for dense generation tasks. Besides, in order to have stability training, the spectral normalization [32] is exploited to regularize the ConvTranspose2d layers of DC-GAN, which controls the weights of modules by the Lipschitz constant.

Concretely, as displayed in Fig. 2, the noise $z$ is input to generator $\mathcal{G}$ to synthesize the alternative samples $\hat{x}$. We first divide the whole DCGAN module into four blocks. Then we plug the attention module at each block to compute the
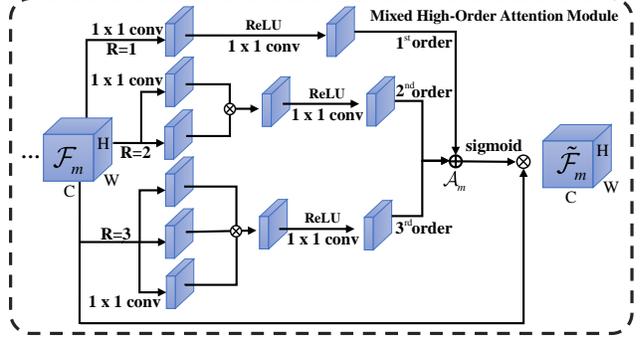


Figure 3. The MHA module of teacher and student in distillation stage.

low-dimensional feature maps $\mathcal{A}_d \in \mathbb{R}^{C/r \times H \times W}$ from original feature maps $\mathcal{F}_g \in \mathbb{R}^{C \times H \times W}$, where $r$ is the scaled scalar, C denotes channel, H and W represent the size of the feature maps. This aims to achieve lightweight feature maps. Next, the encoder is employed to achieve the latent space as follows:

$$\begin{cases} \mathcal{A}_d = \mathrm{Cov}^{1 \times 1}(\mathcal{F}_g) , \\ \Psi = \mathrm{ReLU}(\mathrm{BN}(\mathrm{Cov}^{3 \times 3}(\mathcal{A}_d))) , & (4) \\ \Gamma = \mathrm{ReLU}(\mathrm{BN}(\mathrm{Cov}^{3 \times 3}(\mathrm{MP}(\Psi)))) , \end{cases}$$

where $\Psi$ represents features of intermediate process, $\mathrm{Cov}^{1 \times 1}$ and $\mathrm{Cov}^{3 \times 3}$ denote the convolution with kernel size of $1 \times 1$ and $3 \times 3$, and MP represents maximum pooling.

By Eq. 4, we can get the representation of low-dimensional latent space $\Gamma$ from $\mathcal{A}_d \in \mathbb{R}^{C/r \times H \times W}$, and then decode the space with maximum uppooling (MUP) to achieve spatial-wise attention $\mathcal{A}_s \in \mathbb{R}^{C/r \times H \times W}$. This operation can preserve information of the key locations in the feature to achieve the 2D spatial attention map $\mathcal{A}_s$ as follows:

$$\begin{cases} \Psi = \mathrm{MUP}(\mathrm{ReLU}(\mathrm{BN}(\mathrm{DC}^{3 \times 3}(\Gamma)))) , \\ \mathcal{A}_s = \mathrm{Cov}^{1 \times 1}(\mathrm{ReLU}(\mathrm{BN}(\mathrm{DC}^{3 \times 3}(\Psi)))) , & (5) \end{cases}$$

where $\mathrm{DC}^{3 \times 3}$ denotes the deconvolution with kernel size of $3 \times 3$, MUP represents the maximize unpooling. Then, aggregating the attention maps to the original feature maps to achieve $\tilde{\mathcal{F}}_g$ is formulated as:

$$\tilde{\mathcal{F}}_g = \lambda(\mathrm{Softmax}(\mathcal{A}_s) \times \mathcal{F}_g) + \mathcal{F}_g , \quad (6)$$

where $\lambda$ is the hyperparameter to balance the attention maps with features, which defaults to $5e^{-2}$ in our experiments. More details about the contributions of the attention generator are presented in Tab. 6.

### 4.2. Mixed High-Order Attention Distillation

In the stage of distillation, traditional DFKD methods [5, 11, 30] to solve coarse-grained classification commonly exploit the distribution of output layers due to the

significant inter-class variation (compared to intra-class variation), which enables deep networks to learn generalized discriminatory features of coarse-grained classification. However, the distribution knowledge distillation only exploits category-related information with dark knowledge [16], which lacks semantically relevant information. We argue that this paradigm may not be ideal for FGVC, due to the data-free scenario.

To solve the above difficulties, recent methods commonly exploit attention mechanism [3, 31] to capture the discriminative features of the object. However, the existing FGVC methods of attention mechanism are mainly designed for data-available scenarios, and there is no related research in the data-free scenario. This motivates us to extend this strategy in a data-free setting. Besides, the related attention distillation works [52, 37] only consider the low-order attention information, which only focuses on the local information and cannot capture the complex interactions among parts, resulting in less discriminative attention proposals and failing in capturing the subtle differences among objects. In the data-free scenario, due to the semantic information being sparse [9], we believe that low-order attention distillation cannot fully express the knowledge of the features. Thus we propose to exploit mixed high-order attention (MHA) to distill the aggregated local features and semantic context relation of synthesized FGVC images.

Our mixed high-order attention module is shown in Fig. 3, in which mixed 3-order attention (i.e., $R = 3$) is exploited. The feature $\mathcal{F}_m \in \mathcal{R}^{H \times W \times C}$ is first extracted by three route $1 \times 1$ convolutions to achieve 3-order intermediate representations. In each route, the convolution layer and produced relative representation are the same as the order $R$. Then we multiply the representations of each order to obtain aggregated representations. For each aggregated representation, we exploit RELU and $1 \times 1$ convolution to produce the new map which will be aggregated with global attention maps $\mathcal{A}_m$. At last, the activated global attention map $\mathcal{A}_m$ will be multiplied with the original features $\mathcal{F}_m$ to produce the final attention maps $\tilde{\mathcal{F}}_m = \mathcal{A}_m \times \mathcal{F}_m$.

For teacher and student, their channels may be different. Thus we first exploit the *Adapter* to upgrade the channel of the student to the same number as the teacher, which is also implemented by the $1 \times 1$ 2D convolution. Therefore, at each block of the intermediate layer of $\mathcal{T}$ and $\mathcal{S}$, we exploit mean square error (MSE) to measure the MHAD loss, which is formulated as:

$$\mathcal{L}_{\mathrm{MHAD}} = \frac{1}{N \times C} \sum_{i=1}^{N} \sum_{j=1}^{C} \mathrm{MSE}(\mathcal{F}_m^t, \mathcal{F}_m^s) , \qquad (7)$$

where $N$ and $C$ represent the batch size and channel, while $\mathcal{F}_m^t$ and $\mathcal{F}_m^s$ denote the attention map of an intermediate block of teacher and student, respectively. It should be noted that this strategy is only exploited during our training

**Algorithm 1** The whole pipeline of DFKD-FGVC.
___
**Input**: A pre-trained teacher model $\mathcal{T}$ on real data, generator $\mathcal{G}$ and student network $\mathcal{S}$.
**Output**: A well-trained student network $\mathcal{S}$.

1: **// Ganerator Stage**
2: **for** *number of iterations* **do**
3:    **for** *t steps iterations* **do**
4:       Generate random noise $z \sim \mathcal{N}(0, 1)$ ;
5:       Synthesize supporting sample $\hat{x} = \mathcal{G}(z)$ ;
6:       Optimize the generator by $\mathcal{L}_{\mathrm{BN}}$, and $-\mathcal{L}_{\mathrm{KD}}$;
7:       Freeze $\mathcal{S}$ and $\mathcal{T}$, and update $\mathcal{G}$ by Eq. 10 .
8:    **end for**
9: **end for**
10: **// Distillation Stage**
11: **for** *number of iterations* **do**
12:    **for** *k steps iterations* **do**
13:       Generate random noise $z \sim \mathcal{N}(0, 1)$ ;
14:       Synthesize supporting sample $\hat{x} = \mathcal{G}(z)$ ;
15:       Calculate discrepancy by $\mathcal{L}_{\mathrm{KD}}$, $\mathcal{L}_{\mathrm{MHAD}}$, and $\mathcal{L}_{\mathrm{SFCL}}$.
16:       Freeze $\mathcal{G}$ and $\mathcal{T}$, and update $\mathcal{S}$ by Eq. 11 ;
17:    **end for**
18: **end for**
___

process, which does not participate in the inference. Therefore, this does not affect the efficiency of the model.

### 4.3. Semantic Feature Contrast Learning

Since the pre-trained teacher has a higher discriminative ability than the student, optimizing the student by comparing the features of the teacher is conducive to improving the ability of the student to distinguish right from wrong. Therefore, in our FGVC task, we not only focus on intermediate low-level feature variances but also high-level semantic variances of the penultimate layer. Unlike traditional paradigms [6, 7, 21, 39], which contrast the original [7, 21] and augmentation data or in data-driven scenarios [6, 39]. we exploit high-level semantic features to contrast feature representation of teacher and student and aim to learn the variances between different categories in the data-free scenarios, which are more difficult than data-driven scenarios.

Specifically, in the penultimate layer, we obtain their semantic feature representations and exploit the multi-layer perceptron (MLP) to map the representations to a common space to achieve 2N feature representations as $\mathcal{F}_s = \mathcal{C}(\mathcal{S}(\mathcal{G}(z)))$ and $\mathcal{F}_t = \mathcal{C}(\mathcal{T}(\mathcal{G}(z)))$, where $\mathcal{C}$ is the MLP layer with two hidden linear layers. Then, we normalize the features to a unit hyperspace and measure their similarity as follows:

$$sim(\mathcal{F}_t, \mathcal{F}_s) = \frac{\mathcal{F}_t \cdot \mathcal{F}_s^{\top}}{\|\mathcal{F}_t\| \cdot \|\mathcal{F}_s\|} , \qquad (8)$$

where $\cdot$ denotes the inner (dot) product. The cosine dis-

Table 1. Results of different data-free distillation methods on three fine-grained datasets.

| | Setting | Prior Info. | Compression Info. | | Accuracy | | |
|---|---|---|---|---|---|---|---|
| Method | Data-free | BN | FLOPs | Params. | Aircraft | Cars196 | CUB200 |
| ResNet-34 (T.) | × | × | ~3.67G | ~22M | 70.15 | 84.22 | 76.87 |
| ResNet-18 (S.) | × | × | ~1.82G | ~11M | 68.71 | 77.43 | 58.60 |
| ZSKD [33] | ✓ | × | ~1.82G | ~11M | 37.32 | 26.21 | 30.53 |
| ZSKT [30] | ✓ | × | ~1.82G | ~11M | 51.16 | 28.48 | 31.88 |
| DAFL [5] | ✓ | × | ~1.82G | ~11M | 43.69 | 37.71 | 31.01 |
| DFAD [11] | ✓ | × | ~1.82G | ~11M | 49.51 | 48.72 | 40.15 |
| ADI [47] | ✓ | ✓ | ~1.82G | ~11M | 58.14 | 65.24 | 47.63 |
| DFQ [8] | ✓ | ✓ | ~1.82G | ~11M | 60.22 | 66.14 | 48.43 |
| MAD [10] | ✓ | ✓ | ~1.82G | ~11M | 63.74 | 67.53 | 53.43 |
| CMI [12] | ✓ | ✓ | ~1.82G | ~11M | 63.57 | 68.74 | 53.53 |
| **Ours** | ✓ | ✓ | ~1.82G | ~11M | **65.76** | **71.89** | **56.93** |

tance is used as the similarity metric to measure the relationship between two feature representations for contrastive loss, which is defined as

$$\mathcal{L}_{\text{SFCL}} = \min_{\mathcal{S}} \left\{ -\log \frac{\exp(sim(\mathcal{F}_t^i, \mathcal{F}_s^j)/\tau)}{\sum_k^{2\text{N}} \mathbb{1}_{[k \neq i]} \exp(sim(\mathcal{F}_t^i, \mathcal{F}_s^k)/\tau)} \right\} , \tag{9}$$

where $\mathbb{1}_{[k \neq i]}$ is an indicator function that returns 1 if $i = j$, $i$ and $j \in 2\text{N}$ are the indexes of the samples in the representations, and $\tau$ denotes a temperature parameter. This loss maximizes the representations of the different categories, where the teacher can extract the effective features from noisy images and pull away from the other dissimilar features. Therefore, if one feature of the teacher is viewed as an *anchor*, and the student extracts another representation of this synthetic image as the *positive*. Due to the weak ability of sample representations of the student model, such operation of the student plays a role as augmented images. The other $2(\text{N} - 1)$ features can be viewed as the *negative*. Therefore, the loss $\mathcal{L}_{\text{SFCL}}$ is used to optimize the student to close to the teacher model, i.e., improving the ability of students to distinguish different samples.

### 4.4. Total Objects

In the whole algorithm pipeline 1, we first optimize the generator to synthesize more realistic diverse samples. The total objective of the generator is

$$\min_{\mathcal{G}} \alpha \mathcal{L}_{\text{BN}} - \mathcal{L}_{\text{KD}} . \tag{10}$$

Then, with the above strategy for the generator, we can detail the total objective of the student:

$$\min_{\mathcal{S}} \mathcal{L}_{\text{KD}} + \beta \mathcal{L}_{\text{MHAD}} + \gamma \mathcal{L}_{\text{SFCL}} , \tag{11}$$

where $\alpha$, $\beta$, and $\gamma$ are both hyper-parameters. The training plays an adversarial distillation to optimize both at each iteration.

## 5. Experiments

### 5.1. Datasets and Implementation Details

**Datasets.** To demonstrate the effectiveness of our approach, we conduct experiments on three fine-grained datasets.
**Aircraft.** FGVC-Aircraft [29] contains 100 different aircraft variants formed by 10,000 annotated images, which is divided into two subsets, i.e., the training set with 6,667 images and the testing set with 3,333 images.
**Cars196.** The Stanford Cars dataset [23] contains 16,185 images from 196 categories of cars. The data is split into 8,144 training images and 8,041 testing images.
**CUB200.** The Caltech-UCSD birds dataset (CUB-200-2011) [43] consists of 11,788 annotated images in 200 subordinate categories, including 5,994 images for training and 5,794 images for testing.
**Implementation Details.** Our method is implemented with the PyTorch library. All the models are trained on NVIDIA 3090 GPUs with 24G memory. ResNet-34 [15] is employed as the cumbersome teacher network for all experiments in this paper, and four architectures, i.e., ResNet-18 [15], WRN40-2 [15], MobileNetV2 [36], and ResNet-34 [15] are utilized as students. We first train the generator for 20 steps (i.e., $t$=20) where the generator follows the architecture of DCGAN [35]. Adam [22] is adopted to optimize the generator with an initial learning rate of $1 \times 10^{-3}$ and $beta$ is set 0.5 to 0.99. Then, we train the student 15 steps (i.e., $k$=15) after the generator and optimize the parameters by the SGD optimizer with a momentum of 0.9, a batch size of 64 as default, and a weight decay of $5 \times 10^{-4}$. The initial learning rate starts at $1 \times 10^{-2}$ with cosine annealing for a total of 200 epochs. In the pre-trained stage, due to subtle discrepancies that are difficult to detect, the input images of fine-grained datasets are both resized and randomly cropped to 224×244. In the data-free distillation stage, all the synthetic images are the same as the size of the original input images in the pre-trained stage. As for the
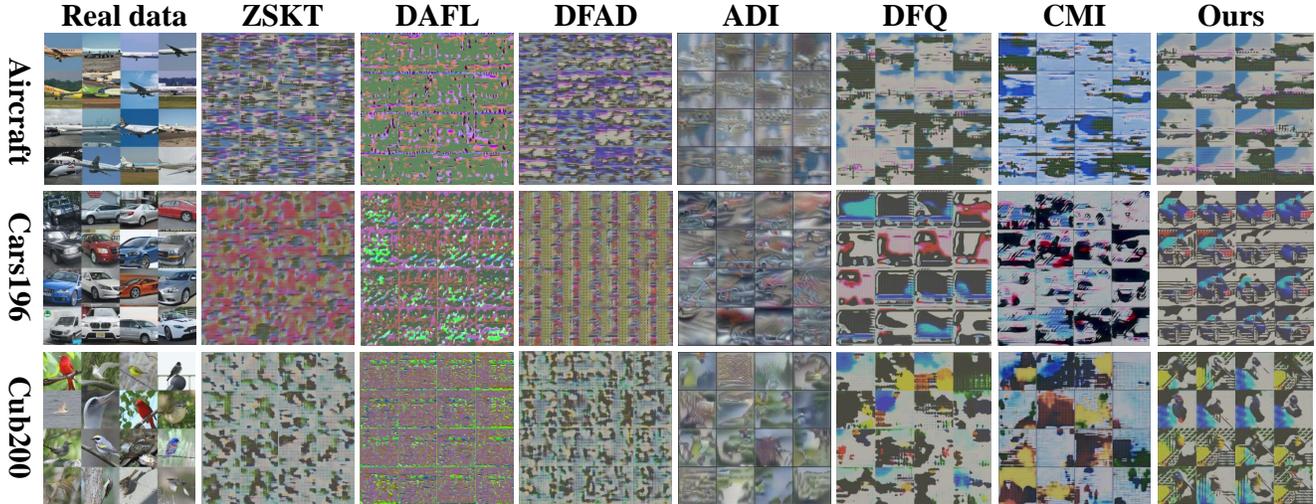
Figure 4. Visualization synthetic images generated by some representative approaches on Aircraft, Cars196, and CUB200 datasets.

hyper-parameters, both $\alpha$, $\beta$, and $\gamma$ are set to 0.3, 10, and 8 by default, respectively. Floating point operations (FLOPs) and parameters (Params) are employed to measure the computation and storage cost of the networks.

## 5.2. Results and Comparisons

As shown in Table 1, we focus on evaluating our approach and other compared methods on three public fine-grained datasets, i.e., Aircraft, Cars196, and Cub200. To evaluate the effectiveness of our proposed method, we conduct fair comparison experiments with two kinds of DFKD methods which are primarily for general classification tasks: (1) Without ($\times$) prior information methods, including ZSKT, DAFL, and DFAD; (2) With ($\checkmark$) prior information methods, including ADI, DFQ, MAD, and CMI. The first two rows of the table show the results of the teacher and student with annotated data supervision in training, which is also our target to achieve by KD. Obviously, the performance of the methods exploiting prior information is better than those without. For example, DFAD only achieves 49.51%, 48.72%, and 40.15% on three datasets, while ADI can achieve 58.14%, 65.24%, and 47.63%. This is mainly because BN regularization has a good performance to inverse and generate relatively realistic images, which is particularly important for downstream distillation. Based on the BN regularization, our approach exploits the spatial attention generator to generate the images with semantic information, which can further improve the performance of the student.

Besides, almost all of the above approaches exploit the vanilla KD (e.g., KL divergence) to transfer the knowledge from the output layer, although they can perform well on coarse-grained classification, but do not perform well on fine-grained classification. Our method mainly adopts two strategies to further improve the performance of the student by about 3% on average, which indicates that vanilla KD alone cannot complete all knowledge transfer, and special design is necessary for FGVC tasks distillation in DFKD. Under identical conditions, thanks to two optimization strategies, i.e., MHAD and SFCL, our approach outperforms the other data-free methods to achieve state-of-the-art performance on three datasets.

Table 2. More comparisons of different architectures' students with ResNet-34 on Aircraft dataset.

| Student | ZSKT | DFAL | DFAD | ADI | DFQ | MAD | CMI | **Ours** |
|---|---|---|---|---|---|---|---|---|
| WRN40-2 | 49.13 | 36.83 | 50.44 | 57.83 | 58.26 | 59.85 | 62.43 | **64.54** |
| MobileNetV2 | 24.39 | 18.51 | 23.01 | 53.66 | 53.93 | 54.61 | 55.04 | **57.37** |
| ResNet-34 | 39.52 | 36.63 | 52.15 | 60.75 | 61.75 | 63.12 | 64.66 | **65.48** |

To verify the generality of our approach, we perform distillation on another three student models with different architectures, including heterogeneous distillation (i.e., WRN40-2 and MobileNetV2) and self-distillation (i.e., Resnet-34). For WRN40-2 and MobileNetV2, we leverage the MLP with two hidden layers to map the dimension to match the teacher and implement our two strategies both in the penultimate layer. As shown in Table 2, our approach can also achieve state-of-the-art performance in different architectures.

## 5.3. Visualization and Analysis

**Synthetic images.** To clearly evaluate the effect of synthesized images, we present a visualization analysis of some representative methods on Aircraft, Cars196, and CUB200 in this section. As we can see from Fig. 4, the first column is the real data for reference. However, for ZSKT, DAFL, and DFAD, there is a big gap between the generated images and the real data. Since ADI, CMI, and Ours both exploit the BN to regularize the features, the synthesized images are more realistic than the other data-free methods, which is beneficial for downstream distillation. With the assis-

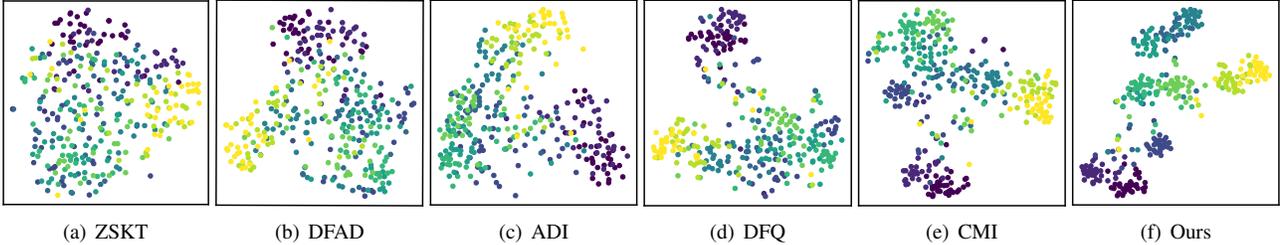(a) ZSKT    (b) DFAD    (c) ADI    (d) DFQ    (e) CMI    (f) Ours

Figure 5. Visualization of t-SNE distribution on Aircraft dataset.

tance of two optimization strategies of MHAD, and SFCL, our approach can generate better and more discriminative foreground images compared to ADI, DFQ, and CMI. For example, we can clearly distinguish the outline of the car and the color of the different areas of the birds.

**t-SNE.** To illustrate the advantages of our approach in synthesizing images having more similar distributions with real images, we sample 10 categories from the Aircraft dataset and visualize the representations of MobileNetV2 by t-SNE as Fig. 5. As shown in Fig. 5(f), our approach gains obviously better representations than the other methods, according to the comparison with each other. Compared the Tab. 1 with Fig. 4, we can conclude that the performance of the student primarily relies on the quality of the synthetic images and the effect of knowledge transfer in DFKD.

**Attention map.** To further verify the effect of our mixed high-order attention (MHA) modules feature selection, we visualized the generated samples through GradCAM [1], as shown in Fig. 6. The first row is the synthesized alternative samples of CUB200 which are generated by our attention module. We can see the fine-grained semantic information of different synthesized birds. For example, we can distinguish different beaks or wings of birds, and different colors of features. When we employ the student embedded with MHA modules to visualize birds' discriminative features by GradCAM, the attention maps are sparse and focus on the discriminative parts, as shown in the second row of the figure. For example, the wings of birds are activated, which indicates that the wings are being paid attention to. We can conclude that MHA modules can focus on contextual semantic information on features which is based on the attention of discriminative features.

### 5.4. Ablation Study

**Contribution of loss.** To verify the contribution of each component, we conduct ablation experiments on the three datasets with ResNet-18, as shown in Table 3. In the first row is the Baseline of each benchmark, which exploits the Eq. 10 to optimize the $\mathcal{G}$, while only optimizing the $\mathcal{S}$ by exploiting the $\mathcal{L}_{\mathrm{KL}}$ to distill the knowledge. Then, adding the $\mathcal{L}_{\mathrm{SFCL}}$ component to the Baseline, the result of each benchmark is improved by 3.07%, 2.33%, 2.92%, respec-
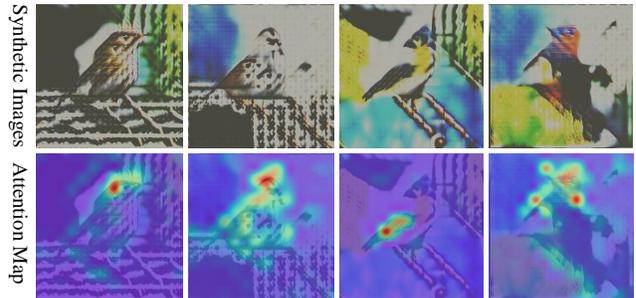


Figure 6. Visualization of synthetic images with attention map generated by GradCAM on CUB200 datasets.
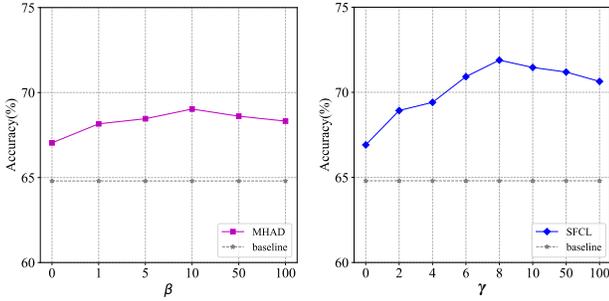
tively. Likewise, when we add $\mathcal{L}_{\mathrm{MHAD}}$ to the baseline, it can also achieve significant improvement. Nevertheless, by comparing both, we can find that the contribution of $\mathcal{L}_{\mathrm{SFCL}}$ is relatively weaker than $\mathcal{L}_{\mathrm{MHAD}}$, which proves the effectiveness of exploiting mixed high-order attention to model discriminative features, which has been ignored by other methods. Finally, we add both to the baseline and obtain the final state-of-the-art effect.

Table 3. The ablation study of our approaches with different components. '+' denotes the add operation.

| Method | Aircraft | Cars196 | Cub200 |
|---|---|---|---|
| Baseline | 60.30 | 64.80 | 51.34 |
| + $\mathcal{L}_{\mathrm{SFCL}}$ | 63.37 | 67.13 | 54.26 |
| + $\mathcal{L}_{\mathrm{MHAD}}$ | 64.86 | 69.92 | 55.71 |
| + $\mathcal{L}_{\mathrm{MHAD}}$ + $\mathcal{L}_{\mathrm{SFCL}}$ | 65.76 | 71.89 | 56.93 |

**Effect of hyper-parameters.** In our optimization, $\alpha$, $\beta$, and $\gamma$ are the major hyper-parameters for balancing the loss terms in our framework. By adjusting the BN hyperparameter in the interval between 0 to 5, we find that the optimal value of $\alpha$ is 0.3. Then, we investigate the effect of $\beta$ and $\gamma$ on the student ResNet-18 on the Cars196 dataset and show the results in Fig. 7. In Fig. 7(a), we set $\gamma$ as 1.0 and vary $\beta$ from 0 to 100, in which 10 is a reasonable parameter verified by our experiments. Then, we set the optimal value of $\beta$ to 10 and vary $\gamma$ from 0 to 100, in which the student network achieves the best performance when $\gamma$ is set to 8, as shown in Fig. 7(b). It is clear that, when using different $\beta$ and $\gamma$, our model stably outperforms the baseline model. The experimental results show that our proposed framework is robust to the different parameters.

---

[1]https://github.com/jacobgil/pytorch-grad-cam.git

(a) $\alpha = 0.3$, $\gamma = 1.0$, adjust $\beta$  (b) $\alpha = 0.3$, $\beta = 10$, adjust $\gamma$

Figure 7. Effect of hyper-parameter $\beta$ and $\gamma$ on Cars196 dataset.

**Control parameter of attention ganerator.** Due to the parameter $\lambda$ being exploited to control the aggregating of attention and feature maps, we perform a group analysis of this parameter. As shown in Tab. 4, we first fix the other parameters, and then the $\lambda$ is set to 0, which indicates that the generator does not exploit the attention mechanism. And the results on the three datasets achieve 63.88, 69.24, and 54.81, respectively. From the interval 0 to $5e^{-2}$, the effect of the generator rises significantly while the effect of the model decays between $5e^{-2}$ and $9e^{-2}$, in which the reasonable parameter is $5e^{-2}$. This indicates that the generator needs to be moderate when employing attention. When the generator pays too much attention to the attention image, it may destroy the original synthesized images resulting in the degradation of the model.

Table 4. The effect of $\lambda$ under different parameters.

| $\lambda$ | Aircraft | Cars196 | CUB200 |
|-----------|----------|---------|--------|
| 0 | 63.88 | 69.24 | 54.81 |
| $1e^{-2}$ | 64.32 | 69.87 | 55.44 |
| $5e^{-2}$ | 65.76 | 71.89 | 56.93 |
| $7e^{-2}$ | 65.02 | 70.95 | 56.30 |
| $9e^{-2}$ | 64.23 | 70.36 | 55.81 |

**Order effectiveness of MHA.** We adopt mixed 3-order attention distillation in our method, which is mainly due to the 3-order attention having the ability to pay attention to the context information. It has more information than the 1-order attention. In this section, we conduct experiments to verify the effect of different orders on different FGVC datasets. As can be seen from Tab 5, when exploiting the 1-order attention distillation, we can only achieve 64.31, 69.26, and 56.12 on three datasets. However, when we exploit the 3-order attention distillation, we can improve the scores of 1.5% on average. What exceeded our expectations is the lower effect when 2-order attention was used. We believe that the 2-order attention mainly focuses on the global information, including the background information, which confuses the foreground attention and reduces the effect of attention.

Table 5. The effect of different orders on different FGVC datasets.

| Order | Aircraft | Cars196 | CUB200 | Avg. |
|-------|----------|---------|--------|------|
| $R = 1$ | 64.31 | 69.26 | 56.12 | 63.23 |
| $R = 2$ | 63.12 | 70.35 | 55.06 | 62.84 |
| $R = 3$ | 65.76 | 71.89 | 56.93 | 64.86 |

### 5.5. Architecture of generator

As illustrated in Fig. 1, the generator with spatial-wise attention modules is adopted in our experiments. Therefore, we detail the architecture of the generator and attention module as indicated in Tab. 6. Concretely, our generator is isomorphic to DCGAN [35]. However, to facilitate the calculation of the spatial-wise attention module, we divide the generator into four blocks. At each block, we exploit spectral normalization to normalize the weights of deconvolution, which aims to stabilize the training of the generator. Then, the encoder-decoder spatial-wise attention module is plugged into each block of the generator, in which the indexes of Maxpool are also used in the MaxUnpool to focus on the key position of synthesized features.

Table 6. The **Left**. Attention Generator Architectures. The noise is mapped to the features which are upsampled to the required image size. The SN denotes the spectral normalization, while SAM represents spatial-wise attention modules corresponding to the **Right**.

| Attention Generator | Spatial-wise Attention Modules |
|---------------------|-------------------------------|
| FC, Reshape, BN | $1 \times 1$ $C \to C/r$ Conv |
| $3 \times 3$, $512 \to 256$, Deconv $\uparrow_{2\times}$, SN, LReLU, SAM | $3 \times 3$, $C/r \to 2C/r$, Conv, BN, ReLU, Maxpool |
| $3 \times 3$, $256 \to 128$, Deconv $\uparrow_{2\times}$, SN, LReLU, SAM | $3 \times 3$, $2C/r \to 4C/r$, Conv, BN, ReLU |
| $3 \times 3$, $128 \to 64$, Deconv $\uparrow_{2\times}$, SN, LReLU, SAM | $3 \times 3$, $4C/r \to 2C/r$, Decov, BN, ReLU, MaxUnpool |
| $3 \times 3$, $64 \to 64$, Deconv $\uparrow_{2\times}$, SN, LReLU, SAM | $3 \times 3$, $2C/r \to C/r$, Decov, BN, ReLU |
| $3 \times 3$, $64 \to 3$, Conv, Tanh | $1 \times 1$, $C/r \to C$, Conv, SoftMax |

\* $C$ is the input channel of each block, while r is scale scalar.

### 6. Conclusion

In this paper, we address the data-free distillation for FGVC. We propose to exploit the generator with spatial attention to synthesize the images with discriminative features. Then, two effective strategies are exploited to optimize the student by MHAD and SFCL, where MHAD captures the discriminative features with context information and SFCL exploits the high-level semantic features to contrast the variances between the different categories. Experimental evidence demonstrates that both approaches can improve the performance of the student on FGVC tasks and outperform other data-free distillation approaches to achieve state-of-the-art performance.

# References

[1] Aharon Ben-Tal, Laurent El Ghaoui, and Arkadi Ne-mirovski. *Robust optimization*, volume 28. Princeton university press, 2009.

[2] Thomas Berg and Peter N Belhumeur. Poof: Part-based one-vs.-one features for fine-grained categorization, face verification, and attribute estimation. In *CVPR*, pages 955–962, 2013.

[3] Sijia Cai, Wangmeng Zuo, and Lei Zhang. Higher-order integration of hierarchical convolutional activations for fine-grained visual categorization. In *ICCV*, pages 511–520, 2017.

[4] Binghui Chen, Weihong Deng, and Jiani Hu. Mixed high-order attention network for person re-identification. In *CVPR*, pages 371–381, 2019.

[5] 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 *ICCV*, pages 3514–3522, 2019.

[6] Liqun Chen, Dong Wang, Zhe Gan, Jingjing Liu, Ricardo Henao, and Lawrence Carin. Wasserstein contrastive representation distillation. In *CVPR*, pages 16296–16305, 2021.

[7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, pages 1597–1607. PMLR, 2020.

[8] Yoojin Choi, Jihwan Choi, Mostafa El-Khamy, and Jungwon Lee. Data-free network quantization with adversarial knowledge distillation. In *CVPR*, pages 710–711, 2020.

[9] Yao Ding, Yanzhao Zhou, Yi Zhu, Qixiang Ye, and Jianbin Jiao. Selective sparse sampling for fine-grained image recognition. In *ICCV*, pages 6599–6608, 2019.

[10] Kien Do, 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. In *Advances in NeurIPS*, 2022.

[11] Gongfan Fang, Jie Song, Chengchao Shen, Xinchao Wang, Da Chen, and Mingli Song. Data-free adversarial distillation. *arXiv preprint arXiv:1912.11006*, 2019.

[12] Gongfan Fang, Jie Song, Xinchao Wang, Chengchao Shen, Xingen Wang, and Mingli Song. Contrastive model inversion for data-free knowledge distillation. In *IJCAI*, 2021.

[13] Weifeng Ge, Xiangru Lin, and Yizhou Yu. Weakly supervised complementary parts models for fine-grained image classification from the bottom up. In *CVPR*, pages 3034–3043, 2019.

[14] Meng-Hao Guo, Tian-Xing Xu, Jiang-Jiang Liu, Zheng-Ning Liu, Peng-Tao Jiang, Tai-Jiang Mu, Song-Hai Zhang, Ralph R Martin, Ming-Ming Cheng, and Shi-Min Hu. Attention mechanisms in computer vision: A survey. *Computational Visual Media*, 8(3):331–368, 2022.

[15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.

[16] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.

[17] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. pages 7132–7141, 2018.

[18] Shaoli Huang, Zhe Xu, Dacheng Tao, and Ya Zhang. Part-stacked cnn for fine-grained visual categorization. In *CVPR*, pages 1173–1182, 2016.

[19] Ruyi Ji, Longyin Wen, Libo Zhang, Dawei Du, Yanjun Wu, Chen Zhao, Xianglong Liu, and Feiyue Huang. Attention convolutional binary neural tree for fine-grained visual categorization. In *CVPR*, pages 10468–10477, 2020.

[20] Zilong Ji, Xiaolong Zou, Xiaohan Lin, Xiao Liu, Tiejun Huang, and Si Wu. An attention-driven two-stage clustering method for unsupervised person re-identification. In *ECCV*, pages 20–36. Springer, 2020.

[21] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in NeurIPS*, 33:18661–18673, 2020.

[22] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2014.

[23] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *ICCV workshops*, pages 554–561, 2013.

[24] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in NeurIPS*, 25, 2012.

[25] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. In *ICLR*, 2016.

[26] Chuanbin Liu, Hongtao Xie, Zheng-Jun Zha, Lingfeng Ma, Lingyun Yu, and Yongdong Zhang. Filtration and distillation: Enhancing region attention for fine-grained visual categorization. In *AAAI*, volume 34, pages 11555–11562, 2020.

[27] Yuang Liu, Wei Zhang, and Jun Wang. Zero-shot adversarial quantization. In *CVPR*, pages 1512–1521, 2021.

[28] Yuang Liu, Wei Zhang, Jun Wang, and Jianyong Wang. Data-free knowledge transfer: A survey. *arXiv preprint arXiv:2112.15278*, 2021.

[29] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.

[30] Paul Micaelli and Amos J Storkey. Zero-shot knowledge transfer via adversarial belief matching. *Advances in NeurIPS*, 32, 2019.

[31] Shaobo Min, Hantao Yao, Hongtao Xie, Zheng-Jun Zha, and Yongdong Zhang. Multi-objective matrix normalization for fine-grained visual recognition. *IEEE Transactions on Image Processing*, 29:4996–5009, 2020.

[32] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. In *ICLR*, 2018.

[33] Gaurav Kumar Nayak, Konda Reddy Mopuri, Vaisakh Shaj, Venkatesh Babu Radhakrishnan, and Anirban Chakraborty. Zero-shot knowledge distillation in deep networks. In *ICML*, pages 4743–4751. PMLR, 2019.

[34] Jongchan Park, Sanghyun Woo, Joon-Young Lee, and In So Kweon. Bam: Bottleneck attention module. *BMVC*, 2018.

[35] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *ICLR*, 2015.

[36] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *CVPR*, pages 4510–4520, 2018.

[37] Changyong Shu, Yifan Liu, Jianfei Gao, Zheng Yan, and Chunhua Shen. Channel-wise knowledge distillation for dense prediction. In *ICCV*, pages 5311–5320, 2021.

[38] Teik Toe Teoh and Zheng Rong. Deep convolutional generative adversarial network. In *Artificial Intelligence with Python*, pages 289–301. Springer, 2022.

[39] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. *arXiv preprint arXiv:1910.10699*, 2019.

[40] Fei Wang, Mengqing Jiang, Chen Qian, Shuo Yang, Cheng Li, Honggang Zhang, Xiaogang Wang, and Xiaoou Tang. Residual attention network for image classification. In *CVPR*, pages 3156–3164, 2017.

[41] Zhuhui Wang, Shijie Wang, Haojie Li, Zhi Dou, and Jianjun Li. Graph-propagation based correlation learning for weakly supervised fine-grained image classification. In *AAAI*, volume 34, pages 12289–12296, 2020.

[42] Xiu-Shen Wei, Yi-Zhe Song, Oisin Mac Aodha, Jianxin Wu, Yuxin Peng, Jinhui Tang, Jian Yang, and Serge Belongie. Fine-grained image analysis with deep learning: A survey. *TPAMI*, 2021.

[43] Peter Welinder, Steve Branson, Takeshi Mita, Catherine Wah, Florian Schroff, Serge Belongie, and Pietro Perona. Caltech-ucsd birds 200. 2010.

[44] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *ECCV*, pages 3–19, 2018.

[45] Di Wu, Chao Wang, Yong Wu, Qi-Cong Wang, and De-Shuang Huang. Attention deep model with multiscale deep supervision for person re-identification. *IEEE Trans. ETCI*, 5(1):70–78, 2021.

[46] Jing Xu, Rui Zhao, Feng Zhu, Huaming Wang, and Wanli Ouyang. Attention-aware compositional network for person re-identification. In *CVPR*, pages 2119–2128, 2018.

[47] 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 *CVPR*, pages 8715–8724, 2020.

[48] Ning Zhang, Jeff Donahue, Ross Girshick, and Trevor Darrell. Part-based r-cnns for fine-grained category detection. In *ECCV*, pages 834–849. Springer, 2014.

[49] Ritchie Zhao, Yuwei Hu, Jordan Dotzel, Chris De Sa, and Zhiru Zhang. Improving neural network quantization without retraining using outlier channel splitting. In *ICML*, pages 7543–7552. PMLR, 2019.

[50] Yifan Zhao, Ke Yan, Feiyue Huang, and Jia Li. Graph-based high-order relation discovery for fine-grained recognition. In *CVPR*, pages 15079–15088, 2021.

[51] Heliang Zheng, Jianlong Fu, Tao Mei, and Jiebo Luo. Learning multi-attention convolutional neural network for fine-grained image recognition. In *ICCV*, pages 5209–5217, 2017.

[52] Zaida Zhou, Chaoran Zhuge, Xinwei Guan, and Wen Liu. Channel distillation: Channel-wise attention for knowledge distillation. *arXiv preprint arXiv:2006.01683*, 2020.

[53] Peiqin Zhuang, Yali Wang, and Yu Qiao. Learning attentive pairwise interaction for fine-grained classification. In *AAAI*, volume 34, pages 13130–13137, 2020.