# Multi-Level Decoupled Relational Distillation for Heterogeneous Architectures

Yaoxin Yang[1]    Peng Ye[1]    Weihao Lin[1]    Kangcong Li[1]    Yan Wen[1]
Jia Hao[1]    Tao Chen[1†]
[1]School of Information Science and Technology, Fudan University

yxyang24@m.fudan.edu.cn, eetchen@fudan.edu.cn

## Abstract

*Heterogeneous distillation is an effective way to transfer knowledge from cross-architecture teacher models to student models. However, existing heterogeneous distillation methods do not take full advantage of the dark knowledge hidden in the teacher's output, limiting their performance. To this end, we propose a novel framework named **M**ulti-**L**evel **D**ecoupled **R**elational **K**nowledge **D**istillation (**MLDR-KD**) to unleash the potential of relational distillation in heterogeneous distillation. Concretely, we first introduce Decoupled Finegrained Relation Alignment (DFRA) in both logit and feature levels to balance the trade-off between distilled dark knowledge and the confidence in the correct category of the heterogeneous teacher model. Then, Multi-Scale Dynamic Fusion (MSDF) module is applied to dynamically fuse the projected logits of multiscale features at different stages in student model, further improving performance of our method in feature level. We verify our method on four architectures (CNNs, Transformers, MLPs and Mambas), two datasets (CIFAR-100 and Tiny-ImageNet). Compared with the best available method, our MLDR-KD improves student model performance with gains of up to 4.86% on CIFAR-100 and 2.78% on Tiny-ImageNet datasets respectively, showing robustness and generality in heterogeneous distillation. Code will be released soon.*

## 1. Introduction

Recently, knowledge distillation (KD) [1], which aims to train a superior lightweight student model by mimicking the teacher model, has been demonstrated to be one of the most effective approaches for model compression [2, 3]. The majority of existing knowledge distillation methods [2, 3, 4, 5, 6] concentrate on the distillation between teacher and student models with homogeneous architectures. However, this narrow focus limits the widespread use of knowledge distillation. On one hand, there continually emerge
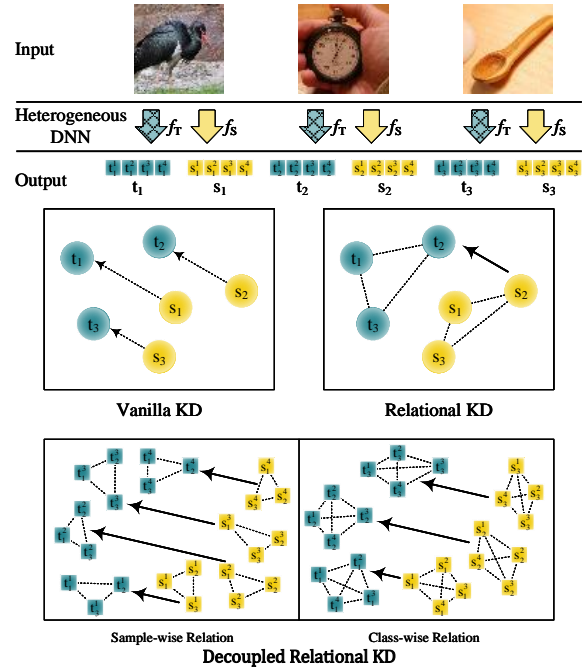


Figure 1. Conceptual comparisons of different knowledge distillation methods. Our Decoupled Relational KD first decouples the logits of teacher and student into multiple finegrained relationships between different classes under each sample and different samples under each class, and then aligns the relationships. In our method, Decoupled Relational KD is applied to both logit and multiscale feature levels (namely MLDR-KD).

new network architectures such as mamba [7]. On the other hand, there exist various pretrained models that have superior performance but different architectures [8, 9, 10]. Consequently, it is essential to explore the potential of knowledge distillation between heterogeneous architectures.

A few recent studies attempt to investigate the feasibility of using heterogeneous teachers for knowledge transfer [13, 14, 15]. Touvron *et al.* [16] achieves successful training of a ViT student model using a CNN teacher model. Ao Wang *et al.* [17] revisits the efficient design of lightweight CNNs from the ViT perspective and emphasizes their promising prospect for mobile devices. Although achieving good re-
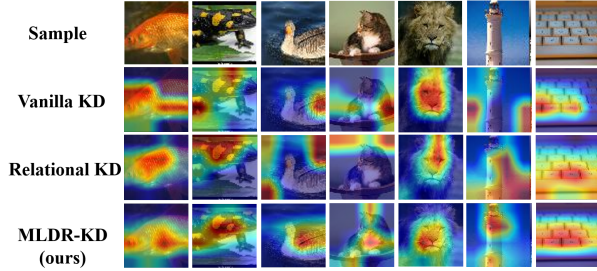
---

†Corresponding authors.

Figure 2. Comparisons of feature visualizations when using kinds of knowledge distillation methods. The teacher is Vision Mamba Tiny [11], the student is ResNet-18 [12]. The direct use of conventional relational KD underperforms on heterogeneous distillation, while our MLDR-KD could greatly improve this problem.

sults, these approaches cannot be extended to various architectures. As a pioneer, Zhiwei Hao *et al.* [18] finds there is a huge gap among feature maps of heterogeneous architecture, resulting in the failure of feature-based knowledge distillation [19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30]. Thus, they propose logit-based generic heterogeneous distillation. Specifically, by increasing the confidence in the correct category of the teacher model, the impact of architectural differences is reduced, and the results are improved. However, this approach somehow weakens the transfer of dark knowledge, which is regarded as very important in knowledge distillation (e.g., whether a sample that is actually a dog or more like a cat), which limits the performance of heterogeneous distillation.

In this paper, we further explore how to effectively transfer the dark knowledge during heterogeneous distillation for the first time. In traditional homogeneous distillation, relational knowledge distillation (RKD) [31] is generally considered as an effective method for transferring dark knowledge, as shown in Fig. 1. RKD aligns correlations or dependencies among multiple instances between the student and teacher networks. However, we find that the direct use of RKD in heterogeneous distillation causes a new problem: the over-amplification of the role of dark knowledge, which may reduce the confidence in the correct category of the teacher model. Since the latter is equally important in heterogeneous distillation due to the variability between architectures, this can directly contribute to the failure of the RKD method, as shown in Fig. 2. Facing such a dilemma, a question naturally arises: *can we effectively transfer the abundant dark knowledge while keeping the confidence of the correct category during heterogeneous distillation?*

To answer this question, we present an innovative framework called Multi-Level Decoupled Relational Knowledge Distillation (MLDR-KD) for heterogeneous distillation. Specifically, we first propose Decoupled Finegrained Relation Alignment (DFRA), in which model logits are first decoupled into multiple finegrained relationships between different categories under each image and different images under each category. Due to the multiple steps finegrained decoupling, the subsequent alignment is sensitive to whether the model classifies correctly, and it can magnify the gap when the classification results of student model and teacher model are not aligned. As a result, our method can well transfer dark knowledge while enhancing the confidence of the classification results during heterogeneous distillation. Further, we apply the DFRA to both logit and feature levels, and present the Multi-Scale Dynamic Fusion (MSDF) module in the feature level. In the MSDF module, the multiscale feature maps of different stages in student model are projected into multiple logits, and a gated network is used to dynamically fuse these logits. As shown in Fig. 2, our method can release the potential of logit-based cross-architectures distillation, where the student model will focus more on information related to the goal.

To illustrate the robustness and generality of our approach, we conduct 12 kinds of experiments between 4 architectures including CNNs, Transformers, MLPs and Mambas. We distill them two by two, with image classification as the evaluation task and acc@1 as the evaluation metric. Compared with the best available method, our MLDR-KD framework improves student model performance with gains of up to 1.43%, 4.86%, 0.93%, 0.83% on CIFAR-100 dataset and 1.57%, 2.78%, 1.61%, 2.13% on Tiny-ImageNet dataset for CNNs, Transformers, MLPs and Mambas architectures under the same conditions, respectively. The ablation study has also demonstrated the effectiveness of our methods. In summary, our main contributions can be summarized as follows:

- We first propose to utilize dark knowledge for heterogeneous distillation. We find that: 1) previous work [18] destroys the dark knowledge present in the teacher model logit, which limits the performance of heterogeneous distillation; 2) The direct use of relational knowledge distillation in traditional homogenous distillation to transfer dark knowledge reduces the confidence in the correct category, bringing about catastrophic performance in heterogeneous distillation.
- To address these, we present a novel framework called Multi-Level Decoupled Relational Knowledge Distillation (MLDR-KD). It consists of Decoupled Finegrained Relation Alignment (DFRA) and Multi-Scale Dynamic Fusion (MSDF) module. Specifically, DFRA enables the student model to learn more finegrained relationships in both logit and feature levels. MSDF module further improves the feature level DFRA by dynamically fusing the predictions of multiscale features of students.
- Extensive experiments across diverse datasets and models consistently verify that MLDR-KD can achieve new state-of-the-art performance. In particular, we extend the MLDR-KD method to the new architecture Mamba, and find our method also performs best, which well illus-

trates the robustness and generality of our method.

## 2. Related work

**Homogeneous Distillation** Hinton *et al.* [1] firstly introduces knowledge distillation to transfer a teacher's knowledge to a student by minimizing their Kullback-Leibler divergence. Following works can be mainly categorized into two pipelines: Feature-based KD and Logits-based KD. To enhance representational capacity, Feature-based KD methods [22, 31] distill knowledge from both intermediate layers and logit outputs. Subsequent works explore various perspectives: CRD [32] emphasizes the structural knowledge of the teacher, while CC [19] identifies instance-level congruent constraints, transferring both instance-level information and inter-instance correlation. Further advancements [28, 33, 34] refine this process with class activation mapping, feature masking, and focal techniques for object detection. Logits-based KD enhances student models by transferring softened targets from teacher models [1, 35]. [36] introduces a Z-score logit standardization method to better capture inter-logit relations to conquer the shared-temperature constraint.

However, in logits-based KD simply using KL divergence is insufficient for exact matching. To tackle the issue, [37] proposes a relation-based loss to preserve inter-class relationships. [38] proposes a novel Cross-Image Relational KD (CIRKD), which focuses on transferring structured pixel-to-pixel and pixel-to-region relations among the whole images. [39] proposes a relational KD framework, Linkless Link Prediction (LLP), to distill knowledge for link prediction with MLPs. These methods seem to solve the problem that dark knowledge is not well transferred in heterogeneous distillation. Nonetheless, these relational distillation methods smooth the logit too much, leading to a reduction in the confidence in the correct category. Moreover, we show that directly transferring conventional relational KD to the heterogeneous distillation setting proves ineffective. Thus, it is a significant necessity to investigate the effective application of relational distillation methods in heterogeneous distillation.

**Heterogeneous Distillation** Heterogeneous distillation allows efficient models to inherit rich representations from powerful teacher models of different architectures, enhancing student model performance and generalization across architectural boundaries. Liu *et al.* [40] pioneers heterogeneous knowledge distillation by aligning the output, attention, and feature spaces of heterogeneous models, assuming identical pixel-level spatial information. To overcome the limitations of this assumption, [41] addresses the architecture gap in cross-architecture distillation by synchronizing the pixel-wise receptive fields of teacher and student networks. However, these methods overlook spatial differences and global context, which FASD [13] addresses

by aligning heterogeneous features and logit mappings between Transformer and Mamba models. However, these approaches do not directly scale to all heterogeneous architectures. Furthermore, OFA-KD [18] explores the feasibility of distilling between multiple architectures. They identified two key limitations in existing methods: lack of latent space alignment, causing inconsistencies in heterogeneous distillation, and absence of adaptive target enhancement, weakening focused knowledge transfer. OFA-KD introduces latent space alignment to eliminate architecture-specific information and adaptive target enhancement to sharpen knowledge transfer, achieving notable gains across diverse models. In this paper, we find that dark knowledge is severely corrupted as OFA-KD changes the distribution of the output logit of the teacher model, which limits the performance of heterogeneous distillation. Therefore, we design a novel heterogeneous relational KD framework called MLDR-KD, which can retain redundant dark knowledge while enhancing confidence in the correct target.

## 3. Methodology

### 3.1. Preliminaries

We start from the original Logit-based Knowledge Distillation (KD) method. Generally, We denote the logit out as $z \in \mathbb{R}^{B \times N}$, where $B$ is the batch size in training and $N$ means the number of categories in dataset. The softmax function is then used to obtain a probability distribution:

$$p_i = \frac{\exp(z_i)}{\sum_{j=1}^{N} \exp(z_j)}, i = 1, 2, \cdots, N \qquad (1)$$

where $p_i$ is the probability distribution of a sample.

In Logit-based KD, the cross-entropy loss $\mathcal{L}_{CE}$ is used to minimize gap between the student model and the ground truth:

$$\mathcal{L}_{CE} = -\sum_{i=1}^{B} \sum_{j=1}^{N} y_{ij} \log(p_{s,ij}) \qquad (2)$$

where $y_{ij}$ is the one-hot encoded true label, and $p_{s,ij}$ is the probability distribution of the student model after softmax.

Student model mimics the teacher model by means of distillation loss $\mathcal{L}_{KL}$. We use the Kullback-Leibler divergence to measure the difference between the student model and the teacher model:

$$\mathcal{D}_{KL}(p_{s,i}||p_{t,i}) = \sum_{j=1}^{N} p_{s,ij} \log \frac{p_{s,ij}}{p_{t,ij}} \qquad (3)$$

$$\mathcal{L}_{KL} = \frac{1}{B} \sum_{i=1}^{B} \mathcal{D}_{KL}(p_{s,i}||p_{t,i}) \qquad (4)$$

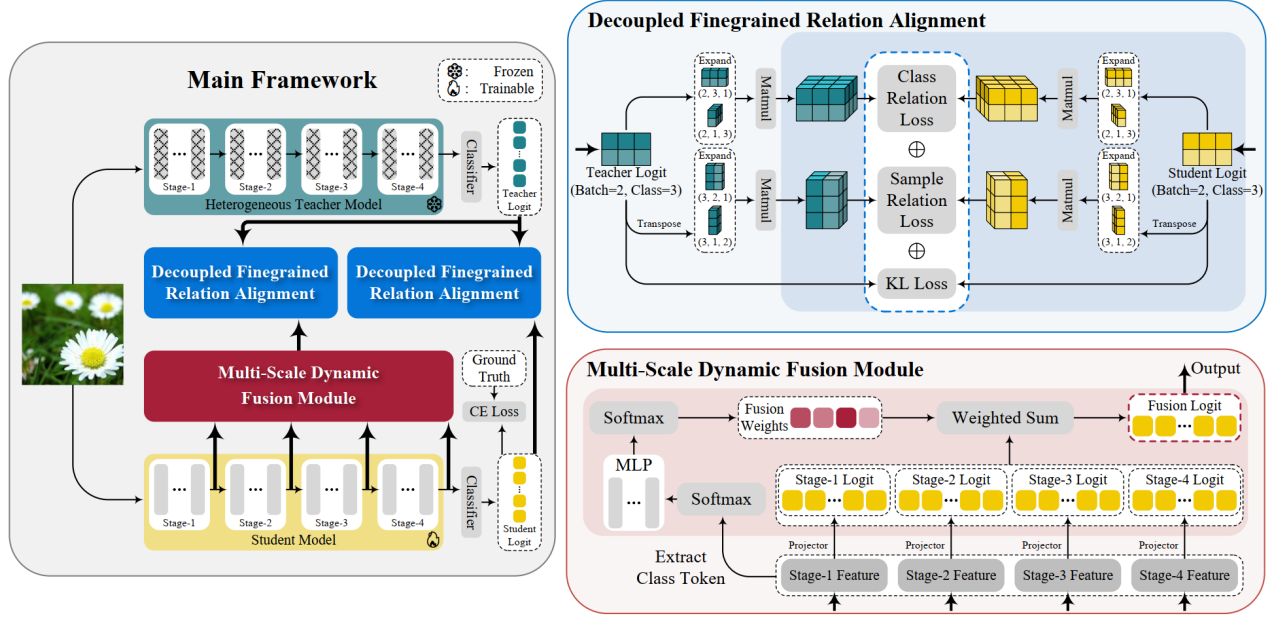where $\mathcal{D}_{KL}(p_{s,i}||p_{t,i})$ is the KL divergence.

Figure 3. Overview of the proposed MLDR-KD framework. It comprises two main components: Decoupled Finegrained Relation Alignment (DFRA), and Multi-Scale Dynamic Fusion (MSDF). In DFRA, after obtaining the logits of teacher and student, we decouple them into class-wise relation and sample-wise relation, and then align these relationships via Kullback-Leibler divergence. DFRA is applied to both logit and feature levels. MSDF further improves the effect of feature-level DFRA by dynamically fusing feature maps of student.

The overall knowledge distillation loss function:

$$\mathcal{L} = \mathcal{L}_{CE} + \lambda \mathcal{L}_{KL} \tag{5}$$

where $\lambda$ is a weighting parameter that balances the cross-entropy loss and distillation loss.

## 3.2. MLDR-KD

The overview of MLDR-KD is depicted in Fig. 3. Our method framework is primarily divided into two modules: the Decoupled Finegrained Relation Alignment (DFRA) and the Multi-Scale Dynamic Fusion (MSDF) Module. Initially, the student model are segmented into multiple stages. After forward inference, the logit outputs of the teacher and the student are obtained. Meanwhile, feature maps of student at each stage are fused by MSDF Module to get a fusion logit. Finally, the fusion logit and the logit output of student will be aligned with logit output of teacher by DFRA in logit and feature levels. We present our two modules of MLDR-KD framework in Sec. 3.2.1 and Sec. 3.2.2.

### 3.2.1 Decoupled Finegrained Relation Alignment

In heterogeneous distillation, it is crucial to balance the correct samples' confidence and the dark knowledge from teacher model. To deal with this problem, we propose DFRA to enhance knowledge transfer between heterogeneous architectures. As shown in Fig. 3, we decouple logit prediction into Class-Wise Relation and Sample-Wise Relation in contrast to exact match. These relations will be aligned in multi level.

**Class-Wise Relation Decoupling** Class-Wise relation represents the degree of similarity among different categories. In this section, we refine this relationship to each sample in the batch to transfer more information (*e.g.* under a particular sample labeled dog, the similarity between cat and elephant). Firstly, we expand logit prediction to three dimensions, which is defined as:

$$\hat{z}_c = \text{Expand}(z/T), \hat{z}_c \in \mathbb{R}^{B \times N \times 1} \tag{6}$$

where $T$ is the soft factor in knowledge distillation. Then we could calculate its self-relation, which is implemented as the scaled product relation:

$$\mathcal{R}_{class} = \text{Softmax}\left(\frac{\hat{z}_c \hat{z}_c^T}{\sqrt{N}}\right), \mathcal{R}_{class} \in \mathbb{R}^{B \times N \times N} \tag{7}$$

where $\mathcal{R}_{class}$ indicates class-wise relation decoupled from initial logit out $z$. $N$ denotes a scaling factor that equals to the number of categories in dataset.

**Sample-Wise Relation Decoupling** The other information then can be decoupled from initial logit out $z$ is sample-wise relation. It's regarded as the degree of similarity between samples under one category (*e.g.* in a batch which of the many samples is more like a dog). Sample-wise relation can be modeled by predictions of a batch of data as follows:

$$\hat{z}_b = \text{Expand}(z^T/T), \hat{z}_b \in \mathbb{R}^{N \times B \times 1} \tag{8}$$

$$\mathcal{R}_{sample} = \text{Softmax}\left(\frac{\hat{z}_b \hat{z}_b^T}{\sqrt{N}}\right), \mathcal{R}_{batch} \in \mathbb{R}^{N \times B \times B} \tag{9}$$
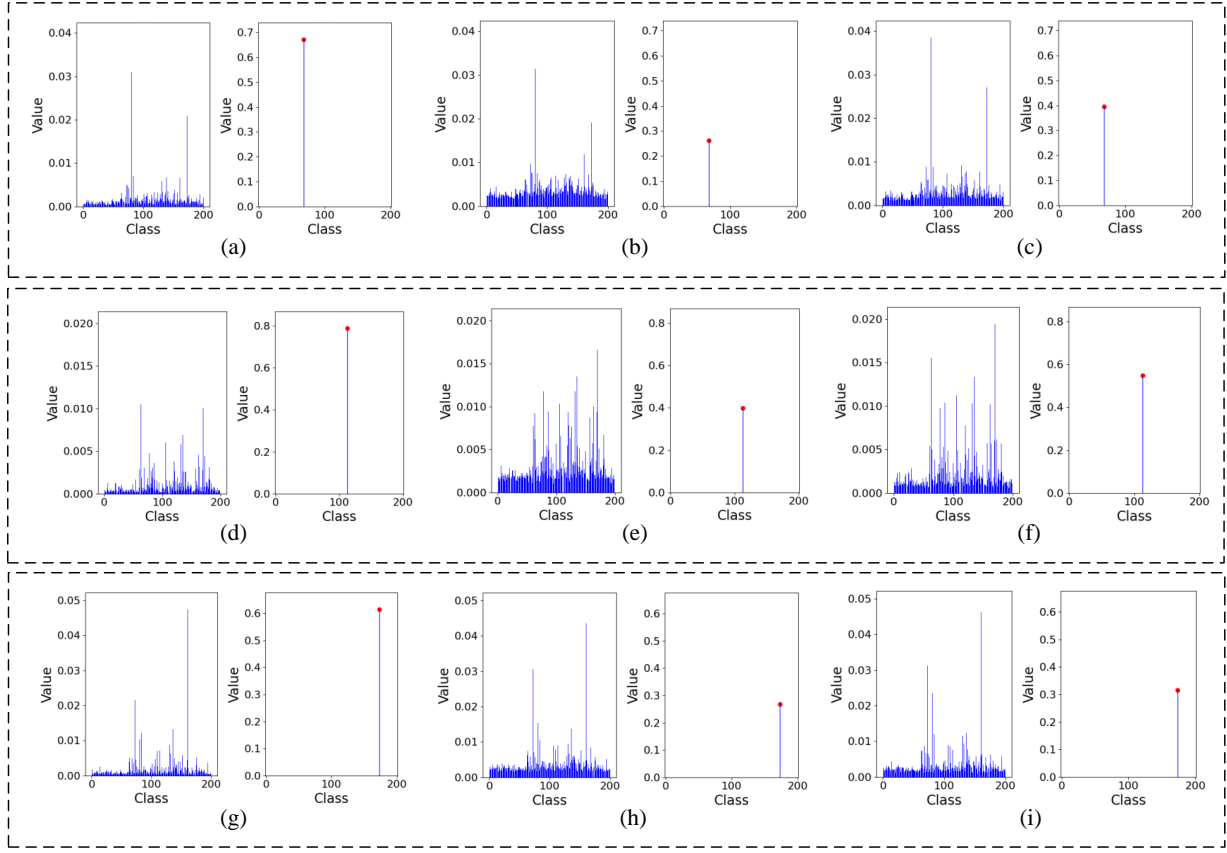
Figure 4. Comparisons of the averaged prediction distribution of all samples of single category among OFA-KD ((a),(d),(g)), RKD ((b),(e),(h)), and our MLDR-KD ((c),(f),(i)). Three black boxes represent three randomly selected categories. In each figure (left), we show the logit of category in addition to the correct category. In each figure (right), the logit of the correct category is displayed. From the figure we can see that our method has high confidence for the correct category while transferring abundant dark knowledge in the teacher model logit.

where $\mathcal{R}_{sample}$ indicates sample-wise relation decoupled from initial logit out $z$.

**Multiple Relation Alignment** Decoupled finegrained relation between heterogeneous student model and teacher model can be aligned by Kullback-Leibler divergence:

$$\mathcal{L}_{class} = \mathcal{L}_{KL}(\mathcal{R}^s_{class}, \mathcal{R}^t_{class}) \tag{10}$$

$$\mathcal{L}_{sample} = \mathcal{L}_{KL}(\mathcal{R}^s_{sample}, \mathcal{R}^t_{sample}) \tag{11}$$

$$\mathcal{L}_{total} = \mathcal{L}_{class} + \mathcal{L}_{sample} + \lambda \mathcal{L}_{KL}(p^s, p^t) \tag{12}$$

where $\lambda$ denotes the balance coefficient. $p^s$ and $p^t$ are the probability distributions of $z^s$ and $z^t$ after softmax function.

Due to the multi-step decoupling of logit, our proposed Decoupled Finegrained Relation Alignment method is robust to enhance the performance of student model. DFRA not only improves the confidence level of the classification results, but also retains a lot of details of the information. Following experiments will further demonstrate the effectiveness of our method.

### 3.2.2 Multi-Scale Dynamic Fusion Module

In heterogeneous distillation, it can transfer more knowledge in addition to the logit level. Specifically, there is a huge gap in the feature maps between heterogeneous models. So it seems feasible to transfer feature-level knowledge in a latent logit space. In other words, feature maps at each stage of student are projected to logit space to be aligned. It can be viewed as training each stage as a separate student. However, the learning abilities of the student models at different stages are disparate because they have different numbers of parameters. It is manifestly inappropriate to assign them the same weighting for learning. So we propose a method that introduce class token in each stage of the student model, denoted as $x_i$ where $i$ is ordinal number of each stage, to dynamically balance weighting of each stage.

The student model is divided into four stages, each of which requires feature matching with the teacher model in logit space, as shown in Figure 2. For each forward inference, the student model outputs features of four stages

Table 1. Results on CIFAR100 dataset. The best results are indicated in bold. For the baseline, most of the experimental results are inherited from OFAKD, while the additional experiments we conducted are marked with *.

| Student Model | From Scratch | Teacher Model | From Scratch | KD [1] | RKD [31] | DKD [35] | OFAKD [18] | **MLDRKD** | Δ |
|---|---|---|---|---|---|---|---|---|---|
| *CNNs-based* | | | | | | | | | |
| ResNet18 | 74.01 | ViT-S | 92.04 | 77.26 | 73.72 | 78.10 | 80.15 | **80.51** | +0.36 |
| | | Swin-T | 89.26 | 78.74 | 74.11 | 80.26 | 80.54 | **81.56** | +1.02 |
| | | Mixer-B/16 | 87.29 | 77.79 | 73.75 | 78.67 | 79.39 | **80.79** | +1.40 |
| | | ViM-S | 87.89* | 78.22* | 77.41* | 79.20* | 79.90* | **80.23** | +0.33 |
| MobileNetV2 | 73.68 | ViT-S | 92.04 | 72.77 | 68.46 | 69.80 | 78.45 | **79.31** | +0.86 |
| | | Mixer-B/16 | 89.26 | 73.33 | 68.95 | 70.20 | 78.78 | **80.21** | +1.43 |
| *Transformers-based* | | | | | | | | | |
| Swin-P | 72.63 | Mixer-B/16 | 87.29 | 75.93 | 69.89 | 76.39 | 78.93 | **80.09** | +1.16 |
| | | ConvNeXt-T | 88.41 | 76.44 | 69.79 | 76.80 | 78.32 | **81.21** | +2.89 |
| | | ViM-S | 87.89* | 78.42* | 72.69* | 79.29* | 79.48* | **79.91** | +0.43 |
| Deit-T | 68.00 | Mixer-B/16 | 87.29 | 71.36 | 70.82 | 73.44 | 73.90 | **78.76** | +4.86 |
| | | ConvNeXt-T | 88.41 | 72.99 | 71.73 | 74.60 | 75.76 | **79.18** | +3.42 |
| | | ViM-S | 87.89* | 73.28* | 70.22* | 74.68* | 76.69* | **77.27** | +0.58 |
| T2t ViT-7 | 74.74 | Mixer-B/16 | 87.29* | 77.43* | 75.76* | 79.53* | 81.54* | **81.61** | +0.07 |
| | | ConvNeXt-T | 88.41* | 79.26* | 75.31* | 79.83* | 82.52* | **82.67** | +0.15 |
| | | ViM-S | 87.89* | 77.39* | 72.53* | 78.48* | 81.38* | **81.47** | +0.09 |
| *MLPs-based* | | | | | | | | | |
| ResMLP-S12 | 66.56 | ConvNeXt-T | 88.41 | 72.25 | 65.82 | 73.22 | 81.22 | **81.96** | +0.74 |
| | | Swin-T | 89.26 | 71.89 | 64.66 | 72.82 | 80.63 | **81.56** | +0.93 |
| | | ViM-S | 87.89* | 80.23* | 78.19* | 80.72* | 80.37* | **80.92** | +0.20 |
| *Mambas-based* | | | | | | | | | |
| ViM-T | 70.99 | ViT-S | 92.04* | 77.55* | 68.85* | 79.58* | 81.24* | **81.87** | +0.63 |
| | | Swin-T | 89.26* | 78.53* | 66.91* | 80.50* | 82.22* | **83.01** | +0.79 |
| | | Mixer-B/16 | 87.29* | 79.34* | 73.08* | 80.57* | 82.19* | **83.02** | +0.83 |
| | | ConvNeXt-T | 88.41* | 80.59* | 66.41* | 82.51* | 82.89* | **82.95** | +0.06 |

$\{f_i\}_{i=1}^4$. We split $\{f_i\}_{i=1}^4$ into the class token $\{x_i\}_{i=1}^4$ for each stage and the architecture-independent feature information $\{\hat{f}_i\}_{i=1}^4$. After that, $\{\hat{f}_i\}_{i=1}^4$ is mapped to the logit space through the projector, denoted as $\{\hat{p}_i\}_{i=1}^4$. We use the global semantic information contained in the class token $\{x_i\}_{i=1}^4$ at each stage to dynamically balance the feature matching under logit space. We apply an MLP layer to generate the balancing weights. MLP layer can be represented as follows:

$$
\begin{aligned}
X_{token} &= \text{Stack}(\{x_i\}_{i=1}^4) \\
X_{hidden} &= \text{GELU}(\text{Linear}(X_{token})) \\
W_{balance} &= \text{Softmax}(\text{Linear}(X_{hidden}))
\end{aligned}
\tag{13}
$$

where $\text{Stack}(\cdot)$ denotes a stacking function for class token aggregation. $\text{GELU}(\cdot)$ indicates an activation function. $\text{Linear}(\cdot)$ is a fully connected layer. In MLP layer, the token sequence $\{x_i\}_{i=1}^4$ is compressed into the vector $X_{token}$. Then with a linear layer and softmax function, we can calculate the balancing weights $W_{balance}$. Further, we use dot product to balance the $\{\hat{p}_i\}_{i=1}^4$, as

$$
\begin{aligned}
P_{stage} &= \text{Stack}(\{\hat{p}_i\}_{i=1}^4) \\
Logit &= W_{balance} \cdot P_{stage}
\end{aligned}
\tag{14}
$$

where $Logit$ is logit output balanced by class token. Finally, we employ DFRA in Sec. 3.2.1 to minimize the gap between $Logit$ and teacher logit $p^t$.

$$
\mathcal{L}_{balance} = \text{DFRA}(Logit, p^t)
\tag{15}
$$

### 3.2.3 Effectiveness Analysis

In Fig. 4, we compare the averaged prediction distribution of all samples of a single category among OFA-KD [18], RKD [31], and our MLDR-KD. Three categories are randomly selected. By comparing each figure left in Fig. 4, we can find that conventional RKD and our MLDR-KD both retain more dark knowledge than the previous heterogeneous distillation method OFA-KD. However, as each figure right shows, conventional RKD reduces the confidence of the student model in the correct category, which leads to its poor performance in heterogeneous distillation. In contrast, the student model trained by our MLDR-KD has high confidence for the correct category while transferring abundant dark knowledge in the teacher model logit, which is consistent with our key observations.

## 4. Experiments

### 4.1. Dataset and Settings

In this section, we will introduce the dataset used in the experiment and the implementation details.

**Datasets** We validate the proposed method on CIFAR-100 [42] and Tiny-ImageNet [43]. The CIFAR-100 dataset comprises 60,000 images divided into 100 categories, with 600 images per category. The size of each image is 32×32. 50,000 images are used as the training set, and 10,000 are used as the test set. The Tiny-ImageNet dataset is a smaller version of the ImageNet dataset. It contains 100,000 images, which are divided into 200 categories. Each category has 500 training images, 50 validation images, and 50 test images. Each image is resized to 64×64.

Table 2. Results on Tiny-ImageNet dataset. The best results are indicated in bold. We conducted all of the additional experiments in baseline. CNN-based experiments are through 100 epochs training. Other experiments are through 300 epochs training.

| Student Model | From Scratch | Teacher Model | From Scratch | KD | OFAKD | MLDRKD | Δ |
|---|---|---|---|---|---|---|---|
| *CNNs-based* | | | | | | | |
| ResNet18 | 63.39 | ViT-S | 80.03 | 65.34 | 65.82 | **67.13** | +1.31 |
| | | Swin-T | 76.13 | 66.20 | 66.94 | **68.51** | +1.57 |
| | | Mixer-B/16 | 69.74 | 64.42 | 65.03 | **66.02** | +0.99 |
| | | ViM-T | 76.13 | 66.69 | 66.62 | **67.61** | +0.92 |
| MobileNetV2 | 63.93 | ViT-S | 80.03 | 66.00 | 65.56 | **66.96** | +0.96 |
| | | Swin-T | 76.13 | 66.51 | 66.60 | **68.06** | +1.46 |
| | | Mixer-B/16 | 69.74 | 64.89 | 65.28 | **65.54** | +0.26 |
| | | ViM-T | 76.13 | 66.26 | 66.14 | **67.24** | +0.98 |
| *Transformers-based* | | | | | | | |
| Swin-P | 65.09 | Mixer-B/16 | 69.74 | 68.67 | 68.24 | **69.10** | +0.43 |
| | | ConvNeXt-T | 72.82 | 66.90 | 67.74 | **68.03** | +0.29 |
| | | ResNet50 | 74.61 | 70.84 | 71.90 | **72.36** | +0.46 |
| | | ViM-T | 76.13 | 70.63 | 70.22 | **70.83** | +0.20 |
| Deit-T | 58.27 | Mixer-B/16 | 69.74 | 64.13 | 68.74 | **69.26** | +0.52 |
| | | ConvNeXt-T | 72.82 | 59.33 | 62.83 | **64.86** | +2.03 |
| | | ResNet50 | 74.61 | 66.72 | 71.89 | **72.29** | +0.4 |
| | | ViM-S | 83.86 | 66.19 | 66.96 | **68.44** | +1.48 |
| | | ViM-T | 76.13 | 67.56 | 68.69 | **71.47** | +2.78 |
| T2t ViT-7 | 64.37 | Mixer-B/16 | 69.74 | 67.34 | 68.85 | **69.36** | +0.51 |
| | | ConvNeXt-T | 72.82 | 65.16 | 66.65 | **69.31** | +2.66 |
| | | ResNet50 | 74.61 | 70.08 | 70.41 | **72.66** | +2.25 |
| | | ViM-T | 76.13 | 69.89 | 70.79 | **72.22** | +1.43 |
| *MLPs-based* | | | | | | | |
| ResMLP-S12 | 65.46 | ConvNeXt-T | 72.82 | 66.37 | 66.74 | **67.23** | +0.49 |
| | | ResNet50 | 74.61 | 72.06 | 70.63 | **73.44** | +1.38 |
| | | ViM-T | 76.13 | 71.58 | 70.31 | **71.72** | +0.14 |
| | | Swin-T | 76.13 | 71.70 | 73.09 | **73.21** | +0.12 |
| | | ViT-S | 80.03 | 70.32 | 69.64 | **71.93** | +1.61 |
| *Mambas-based* | | | | | | | |
| ViM-T | 61.85 | ViT-S | 80.03 | 67.66 | 72.84 | **74.97** | +2.13 |
| | | Swin-T | 76.13 | 70.53 | 72.08 | **73.31** | +1.23 |
| | | Mixer-B/16 | 69.74 | 65.59 | 69.63 | **70.55** | +0.92 |

**Implementation Details** To validate the generality of our method, we conduct experiments with different student and teacher models. For student models, CNN-based ResNet18 [12], MobileNet-v2 [44], Transformers-based Swin-p [45], Deit-t [46], T2t Vit-7 [47], MLP-based ResMLP-S12 [48], and Mamba-based Vim-t [11], are selected. For teacher models, Resnet50 [12], Vit-S [49], Swin-T [45], Mixer-B/16 [50], and ConvNeXt-T [51], are considered. For CNNs, the SGD is adopted as the optimizer, with a base learning rate of 0.05. For Transformers, MLPs, and Mambas, the Adamw [52] is adopted as the optimizer, with a base learning rate of 5e-4. The cosine learning rate decay strategy is used. For all datasets, we set the batch size as 128. The training epoch number of CIFAR-100 is 300 for all models. For Tiny-ImageNet, CNNs are trained with 100 epochs, whereas ViTs, MLPs, and Mambas are trained with 300 epochs. All experiments are conducted using Nvidia RTX 3090 GPU.

## 4.2. Results and Analysis

**Results on CIFAR-100** We first conduct experiments on the CIFAR-100 dataset. Comparisons with the baselines are presented in Table 1. It can be observed that our method can improve the performance of commonly used CNNs-based student models by 0.33% to 1.43%. Moreover, our method achieves remarkable results on Transformers-based student
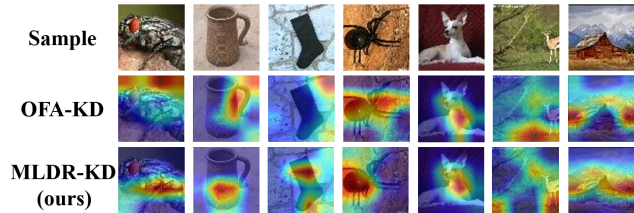


Figure 5. Comparisons of feature visualizations between OFA-KD and our MLDR-KD. The teacher is Vision Mamba Tiny, the student is ResNet-18. Clearly, our approach makes the student model more focused on the target across various samples.

models, especially on the Mixer-B/16 and Deit-t pair, where the accuracy is raised by 4.86%. Compared with KD, DKD, RKD and OFAKD, our method's improvement ranges from 0.09% to 4.86%. For the less prevalent MLPs-based student models, our method also achieves an improvement in accuracy by 0.20% to 0.93% compared with OFAKD. Moreover, the scarcely explored Mambas-based student models are significantly improved by our method, further verifying its effectiveness and generality. Overall, our proposed method achieves state-of-the-art performance on all different student architectures.

**Results on Tiny-ImageNet** To assess the capability of our approach in coping with larger datasets, we expand experiments to the Tiny-ImageNet dataset. To align with the CIFAR-100 dataset, we select corresponding teacher

Table 3. Impact of the number of stages in MSDF (Multi-Scale Dynamic Fusion).

| Number of stage | ACC@1 |
|---|---|
| 0 | 65.33 |
| 1 | 66.03 |
| 2 | 66.94 |
| 3 | 67.08 |
| 4 | **67.13** |

Table 4. Effect of CWRD (Class-Wise Relation Decoupling) and SWRD (Sample-Wise Relation Decoupling) in DFRA (Decoupled Finegrained Relation Alignment).

| MSDF Module | CWRD | SWRD | Acc@1 |
|---|---|---|---|
| ✓ | ✗ | ✗ | 66.76 |
| ✓ | ✗ | ✓ | 66.98 |
| ✓ | ✓ | ✗ | 66.93 |
| ✓ | ✓ | ✓ | **67.13** |

Table 5. More ablation studies in feature and logit levels.

| Feature level | Logit level | MSDF | DFRA | Acc@1 |
|---|---|---|---|---|
| ✓ | ✗ | ✗ | ✗ | 65.98 |
| ✓ | ✗ | ✓ | ✗ | 66.43 |
| ✓ | ✗ | ✗ | ✓ | 66.60 |
| ✓ | ✗ | ✓ | ✓ | 66.96 |
| ✗ | ✓ | ✗ | ✗ | 65.23 |
| ✗ | ✓ | ✗ | ✓ | 65.33 |
| ✓ | ✓ | ✗ | ✗ | 66.68 |
| ✓ | ✓ | ✓ | ✗ | 66.76 |
| ✓ | ✓ | ✗ | ✓ | 66.91 |
| ✓ | ✓ | ✓ | ✓ | **67.13** |

and student models from CNNs-based, Transformers-based, MLPs-based, and Mambas-based models. As no previous baseline results are available, we compare our method with the baselines KD and OFAKD by reproducing them on Tiny-Imagenet. The results are showed in Table 2.

Our method exhibits more stable accuracy improvements than the baselines from the results on the Tiny-ImageNet dataset. Specifically, the accuracy improvement ranges from 0.26% to 1.57% on the CNNs-based student models and 0.20% to 2.78% on the Transformers-based student models, especially achieving an improvement of 2.78% in the architecture pair ViM-T-to-Deit-T. Additionally, there is a significant improvement in our newly added student models T2t ViT-7 compared to the baseline methods. On the MLPs-based student models, our method could achieve an accuracy improvement of 0.12% to 1.61% compared to the baseline methods, particularly attaining the highest improvement of 1.61% on the ResNet50 teacher model. For the latest Mamba-based student models, our method still presents considerable accuracy improvements.

Compared with the results on the CIFAR-100 dataset, the accuracy on the Tiny-ImageNet dataset exhibits advanced stability, indicating the advantages of our fine-grained design over traditional methods when dealing with larger datasets. Moreover, it can be observed that our method has more obvious improvements when applied to larger and more complicated models (such as ConvNeXt-T and ResNet50), validating that our method is potentially practical in further boosting off-the-shelf high-performance models, which are generally large and complicated. Similar to the case on the CIFAR-100 dataset, our method is applicable to the Mamba-based student models, further illustrating the generalization ability of the proposed heterogeneous distillation method.

**Visualization** A visual comparison between our method and baselines on the Tiny-Imagenet dataset is illustrated in Fig. 5. The students trained via our MLDR-KD can better learn from heterogeneous teachers. For example, even though the target occupies a small portion of the sample image, our student model is always able to focus on the information related to the target. The student model's attention does not diverge to distracting information. This is a strong indication that the student model, after our MLDR-KD, is well able to assimilate knowledge from heterogeneous teachers.

### 4.3. Ablation study

Ablative experiments are designed to verify the effectiveness of the proposed MLDR-KD, shown in Table. 3, Table. 4 and Table. 5. In this part, all experiments are conducted on Tiny-Imagenet, with ViT-S as the teacher model, Resnet18 as the student model.

**Number of stages of student** In Table. 3, we conduct experiments to explore the impact of stages in student model. In order to accommodate different architectures, we divide the student model into a maximum of 4 stages. We chose stages 0 to 4 to compare the difference. We find that as the number of stages increases, the improvement effect of our methodology enhances. Four stages are the most potent.

**Validity of CWRD and SWRD in DFRA.** In order to verify the validity of proposed DFRA, we ablate our method in the presence of both feature and logit levels, shown in Table. 4. From the results, CWRD and SWRD have improved by 0.22% and 0.17% relative to the original, respectively. Both of them can enhance 0.37% in performance. Evidently, CWRD and SWRD in DFRA have an indispensable role to play.

**Effect of our MLDR-KD in feature or logit level** Our MLDR-KD improves heterogeneous distillation in both feature and logit levels. We study this improvement in this part. In Table. 5, We ablate different modules at different levels. In a side-by-side comparison (*e.g.* line 1-4), both of our methods MSDF and DFRA are effective when applied to only one level or to both levels. Vertical comparisons (*e.g.* line 4 and 10) show that applying our methodology to multiple levels is the most effective.

## 5. Conclusion

In this paper, we propose Multi-Level Decoupled Relational Knowledge Distillation (MLDR-KD), a novel approach to balance the trade-off between dark knowledge and the confidence in the correct category of the teacher model for heterogeneous architectures. Specifically, DFRA is designed to align finegrained relationship for heterogeneous architectures in feature and logit level. The MSDF module is fur-

ther introduced to improve DFRA performance by fusing feature maps of student in feature level. Extensive experiments show the robustness and generality of our MLDR-KD. Our future work involves how to efficiently take full advantage of feature information to further enhance the proposed MLDR-KD.

# References

[1] Geoffrey Hinton. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 1, 3, 6

[2] Yehui Tang, Kai Han, Yunhe Wang, Chang Xu, Jianyuan Guo, Chao Xu, and Dacheng Tao. Patch slimming for efficient vision transformers, 2022. 1

[3] Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be your own teacher: Improve the performance of convolutional neural networks via self distillation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3713–3722, 2019. 1

[4] Mary Phuong and Christoph H Lampert. Distillation-based training for multi-exit architectures. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1355–1364, 2019. 1

[5] Yunteng Luan, Hanyu Zhao, Zhi Yang, and Yafei Dai. Msd: Multi-self-distillation learning via multi-classifiers within deep neural networks. *arXiv preprint arXiv:1911.09418*, 2019. 1

[6] Xiatian Zhu, Shaogang Gong, et al. Knowledge distillation by on-the-fly native ensemble. *Advances in neural information processing systems*, 31, 2018. 1

[7] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023. 1

[8] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2022. 1

[9] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. Simmim: A simple framework for masked image modeling. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9653–9663, 2022. 1

[10] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. BEiT: BERT pre-training of image transformers. In *International Conference on Learning Representations*, 2022. 1

[11] Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang. Vision mamba: Efficient visual representation learning with bidirectional state space model. In *Forty-first International Conference on Machine Learning*, 2024. 2, 7

[12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 2, 7

[13] Rui Yu, Runkai Zhao, Jiagen Li, Qingsong Zhao, Songhao Zhu, HuaiCheng Yan, and Meng Wang. Unleashing the potential of mamba: Boosting a lidar 3d sparse detector by using cross-model knowledge distillation. *arXiv preprint arXiv:2409.11018*, 2024. 1, 3

[14] Jiaheng Liu, Chenchen Zhang, Jinyang Guo, Yuanxing Zhang, Haoran Que, Ken Deng, Zhiqi Bai, Jie Liu, Ge Zhang, Jiakai Wang, Yanan Wu, Congnan Liu, Wenbo Su, Jiamang Wang, Lin Qu, and Bo Zheng. Ddk: Distilling domain knowledge for efficient large language models. *ArXiv*, abs/2407.16154, 2024. 1

[15] Junxiong Wang, Daniele Paliotta, Avner May, Alexander M Rush, and Tri Dao. The mamba in the llama: Distilling and accelerating hybrid models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 1

[16] Hugo Touvron, Matthieu Cord, and Hervé Jégou. Deit iii: Revenge of the vit. In *European conference on computer vision*, pages 516–533. Springer, 2022. 1

[17] Ao Wang, Hui Chen, Zijia Lin, Jungong Han, and Guiguang Ding. Repvit: Revisiting mobile cnn from vit perspective. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15909–15920, 2024. 1

[18] Zhiwei Hao, Jianyuan Guo, Kai Han, Yehui Tang, Han Hu, Yunhe Wang, and Chang Xu. One-for-all: Bridge the gap between heterogeneous architectures in knowledge distillation. *Advances in Neural Information Processing Systems*, 36, 2024. 2, 3, 6

[19] Baoyun Peng, Xiao Jin, Jiaheng Liu, Dongsheng Li, Yichao Wu, Yu Liu, Shunfeng Zhou, and Zhaoning Zhang. Correlation congruence for knowledge distillation. In *Proceedings of the IEEE/CVF International*

*Conference on Computer Vision*, pages 5007–5016, 2019. 2, 3

[20] Byeongho Heo, Jeesoo Kim, Sangdoo Yun, Hyojin Park, Nojun Kwak, and Jin Young Choi. A comprehensive overhaul of feature distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 2

[21] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *International Conference on Learning Representations*, 2022. 2

[22] 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. 2, 3

[23] Byeongho Heo, Minsik Lee, Sangdoo Yun, and Jin Young Choi. Knowledge transfer via distillation of activation boundaries formed by hidden neurons. In *AAAI Conference on Artificial Intelligence*, 2018. 2

[24] Junho Yim, Donggyu Joo, Ji-Hoon Bae, and Junmo Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7130–7138, 2017. 2

[25] Sungsoo Ahn, Shell Xu Hu, Andreas C. Damianou, Neil D. Lawrence, and Zhenwen Dai. Variational information distillation for knowledge transfer. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9155–9163, 2019. 2

[26] Defang Chen, Jianhan Mei, Hailin Zhang, C. Wang, Yan Feng, and Chun Chen. Knowledge distillation with the reused teacher classifier. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11923–11932, 2022. 2

[27] Pengguang Chen, Shu Liu, Hengshuang Zhao, and Jiaya Jia. Distilling knowledge via knowledge review. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5006–5015, 2021. 2

[28] Ziyao Guo, Haonan Yan, Hui Li, and Xiaodong Lin. Class attention transfer based knowledge distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11868–11877, 2023. 2, 3

[29] Zheng Li, Jingwen Ye, Mingli Song, Ying Huang, and Zhigeng Pan. Online knowledge distillation for

efficient pose estimation. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11720–11730, 2021. 2

[30] Sihao Lin, Hongwei Xie, Bing Wang, Kaicheng Yu, Xiaojun Chang, Xiaodan Liang, and G. Wang. Knowledge distillation via the target-aware transformer. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10905–10914, 2022. 2

[31] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3967–3976, 2019. 2, 3, 6

[32] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. *arXiv preprint arXiv:1910.10699*, 2019. 3

[33] Zhendong Yang, Zhe Li, Mingqi Shao, Dachuan Shi, Zehuan Yuan, and Chun Yuan. Masked generative distillation. In *European Conference on Computer Vision*, pages 53–69. Springer, 2022. 3

[34] Zhendong Yang, Zhe Li, Xiaohu Jiang, Yuan Gong, Zehuan Yuan, Danpei Zhao, and Chun Yuan. Focal and global knowledge distillation for detectors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4643–4652, 2022. 3

[35] 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*, pages 11953–11962, 2022. 3, 6

[36] Shangquan Sun, Wenqi Ren, Jingzhi Li, Rui Wang, and Xiaochun Cao. Logit standardization in knowledge distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15731–15740, 2024. 3

[37] Tao Huang, Shan You, Fei Wang, Chen Qian, and Chang Xu. Knowledge distillation from a stronger teacher. *Advances in Neural Information Processing Systems*, 35:33716–33727, 2022. 3

[38] Chuanguang Yang, Helong Zhou, Zhulin An, Xue Jiang, Yongjun Xu, and Qian Zhang. Cross-image relational knowledge distillation for semantic segmentation, 2022. 3

[39] Zhichun Guo, William Shiao, Shichang Zhang, Yozen Liu, Nitesh V. Chawla, Neil Shah, and Tong Zhao. Linkless link prediction via relational distillation, 2023. 3

[40] Yufan Liu, Jiajiong Cao, Bing Li, Weiming Hu, Jingting Ding, and Liang Li. Cross-architecture knowledge distillation. In *Proceedings of the Asian conference on computer vision*, pages 3396–3411, 2022. 3

[41] Weisong Zhao, Xiangyu Zhu, Zhixiang He, Xiao-Yu Zhang, and Zhen Lei. Cross-architecture distillation for face recognition. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 8076–8085, 2023. 3

[42] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 6

[43] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015. 6

[44] Andrew Howard, Andrey Zhmoginov, Liang-Chieh Chen, Mark Sandler, and Menglong Zhu. Inverted residuals and linear bottlenecks: Mobile networks for classification, detection and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018. 7

[45] 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*, pages 10012–10022, 2021. 7

[46] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, pages 10347–10357. PMLR, 2021. 7

[47] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 558–567, 2021. 7

[48] Hugo Touvron, Piotr Bojanowski, Mathilde Caron, Matthieu Cord, Alaaeldin El-Nouby, Edouard Grave, Gautier Izacard, Armand Joulin, Gabriel Synnaeve, Jakob Verbeek, et al. Resmlp: Feedforward networks for image classification with data-efficient training. *IEEE transactions on pattern analysis and machine intelligence*, 45(4):5314–5321, 2022. 7

[49] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020. 7

[50] Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. Mlp-mixer: An all-mlp architecture for vision. *Advances in neural information processing systems*, 34:24261–24272, 2021. 7

[51] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11976–11986, 2022. 7

[52] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2017. 7