# CrossKD: Cross-Head Knowledge Distillation for Object Detection

Jiabao Wang[1*], Yuming Chen[1*], Zhaohui Zheng[1], Xiang Li[2,1], Ming-Ming Cheng[2,1], Qibin Hou[2,1†]

[1]VCIP, College of Computer Science, Nankai University
[2]NKIARI, Shenzhen Futian

https://github.com/jbwang1997/CrossKD

## Abstract

*Knowledge Distillation (KD) has been validated as an effective model compression technique for learning compact object detectors. Existing state-of-the-art KD methods for object detection are mostly based on feature imitation. In this paper, we present a general and effective prediction mimicking distillation scheme, called CrossKD, which delivers the intermediate features of the student's detection head to the teacher's detection head. The resulting cross-head predictions are then forced to mimic the teacher's predictions. This manner relieves the student's head from receiving contradictory supervision signals from the annotations and the teacher's predictions, greatly improving the student's detection performance. Moreover, as mimicking the teacher's predictions is the target of KD, CrossKD offers more task-oriented information in contrast with feature imitation. On MS COCO, with only prediction mimicking losses applied, our CrossKD boosts the average precision of GFL ResNet-50 with 1× training schedule from 40.2 to 43.7, outperforming all existing KD methods. In addition, our method also works well when distilling detectors with heterogeneous backbones.*

## 1. Introduction

Knowledge Distillation (KD), serving as a model compression technique, has been deeply studied in object detection [5, 13, 29, 31, 56, 60, 61, 74, 75] and has received excellent performance recently. According to the distillation position of the detectors, existing KD methods can be roughly classified into two categories: prediction mimicking and feature imitation. Prediction mimicking (See Fig. 1(a)) was first proposed in [24], which points out that the smooth distribution of the teacher's predictions is more comfortable for the student to learn than the Dirac distribu-
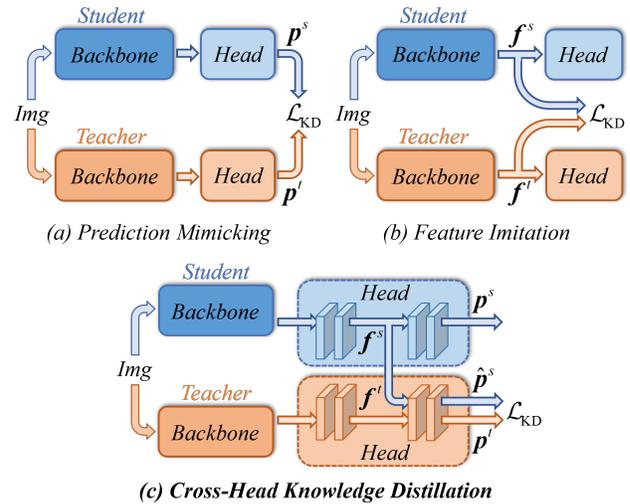


Figure 1. Comparisons between conventional KD methods and our CrossKD. Rather than explicitly enforcing the consistency between the intermediate feature maps or the predictions of the teacher-student pair, CrossKD implicitly builds the connection between the heads of the teacher-student pair to improve the distillation efficiency.

tion of the ground truths. In other words, prediction mimicking forces the student to resemble the prediction distribution of the teacher. Differently, feature imitation (See Fig. 1(b)) follows the idea proposed in FitNet [54], which argues that intermediate features contain more information than the predictions from the teacher. It aims to enforce the feature consistency between the teacher-student pair.

Prediction mimicking plays a vital role in distilling object detection models. However, it has been observed to be inefficient than feature imitation for a long time. Recently, Zheng et al. [74] proposed a localization distillation (LD) method that improves prediction mimicking by transferring localization knowledge, which pushes the prediction mimicking to a new level. Despite just catching up with the advanced feature imitation methods, e.g., PKD [5], LD shows that prediction mimicking has the ability to transfer task-

---

*(a) ground-truth*     *(b) teacher predictions*     *(c) prediction mimicking*     *(d) CrossKD*
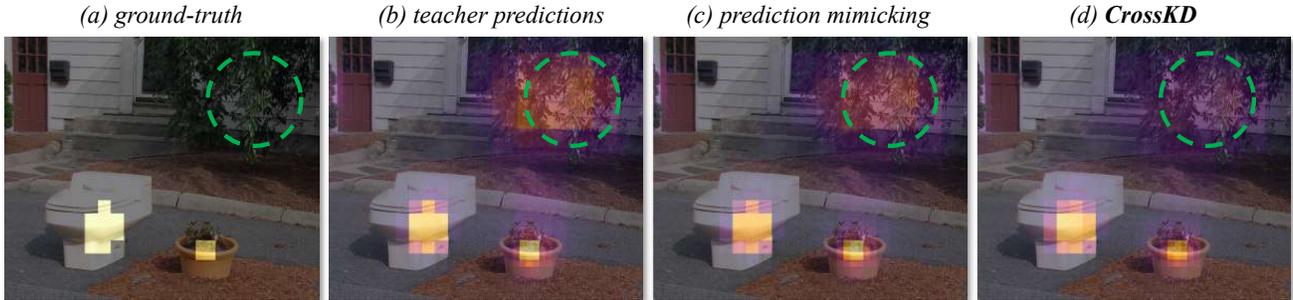
Figure 2. Visualizations of the classification predictions from the GFL [35]. (a) and (b) are ground truth and distillation targets. (c) and (d) are the classification outputs predicted by models training with conventional prediction mimicking and proposed CrossKD. In the green circled areas, the distillation targets predicted by the teacher have a large discrepancy with the ground-truth targets assigned to the student. prediction mimicking forces the student to mimic the teacher, while CrossKD can smooth the mimicking process.

specific knowledge, which benefits the student from the orthogonal aspect to feature imitation. This motivates us to further explore and improve prediction mimicking.

Through investigation, we observe that conventional prediction mimicking may suffer from a conflict between the ground-truth targets from the student's assigner and the distillation targets predicted from the teacher. When training a detector with prediction mimicking, the student's predictions are forced to mimic both the ground-truth targets and the teacher's predictions simultaneously. However, the distillation targets predicted by the teacher usually have a large discrepancy with the ground-truth targets assigned to the student. As shown in Fig. 2(a) and Fig. 2(b), the teacher produces class probabilities in the green circled areas, which conflicts with the ground-truth targets assigned to the student. As a result, the student detector experiences a contradictory learning process during distillation, which seriously interferes with optimization.

To alleviate the above conflict, previous prediction mimicking methods [13, 19, 74] tend to conduct the distillation within regions containing mediate teacher-student discrepancies. However, we argue that the heavily uncertain regions generally accommodate more information that is beneficial to the student. In this paper, we present a novel cross-head knowledge distillation pipeline, abbreviated as *CrossKD*. As illustrated in Fig. 1(c), We propose to feed the intermediate features from the head of the student to that of the teacher, yielding the cross-head predictions. Then, the KD operations can be conducted between the new cross-head predictions and the teacher's predictions.

Despite its simplicity, CrossKD offers the following two main advantages. First, since both the cross-head predictions and the teacher's predictions are produced by sharing part of the teacher's detection head, the cross-head predictions are relatively consistent with the teacher's predictions. This relieves the discrepancy between the teacher-student pair and enhances the training stability of prediction mimicking. In addition, as mimicking the teacher's predictions

is the target of KD, CrossKD is theoretically optimal and can offer more task-oriented information compared with feature imitation. Both advantages enable our CrossKD to efficiently distill knowledge from the teacher's predictions and hence result in even better performance than previous state-of-the-art feature imitation methods.

Without bells and whistles, our method can significantly boost the performance of the student detector, achieving a faster training convergence. Comprehensive experiments on the COCO [40] dataset are conducted in this paper to elaborate the effectiveness of CrossKD. Specifically, with only prediction mimicking losses applied, CrossKD achieves 43.7 AP on GFL with $1\times$ training schedule, which is 3.5 AP higher than the baseline, surpassing all previous state-of-the-art object detection KD methods. Moreover, experiments also indicate our CrossKD is orthogonal to feature imitation methods. By combining CrossKD with the state-of-the-art feature imitation method, like PKD [5], we further achieve 43.9 AP on GFL. Furthermore, we also show that our method can be used to distill detectors with heterogeneous backbones and performs better than other methods.

## 2. Related Work

### 2.1. Object Detection

Object detection is one of the most fundamental computer vision tasks, which requires recognizing and localizing objects simultaneously. Modern object detectors can be briefly divided into two categories: one-stage [3, 10, 11, 35, 39, 51, 57, 70] detectors and two-stage [8, 17, 18, 20, 21, 38, 52, 58, 73] detectors. Among them, one-stage detectors, also known as dense detectors, have emerged as the mainstream trend in detection due to their excellent speed-accuracy trade-off.

Dense object detectors have received great attention since YOLOv1 [49]. Typically, YOLO series detectors [2, 16, 44, 49–51] attempt to balance the model size and their accuracy to meet the requirement of real-world applications. Anchor-free detectors [27, 57, 76] attempt to dis-

card the design of anchor boxes to avoid time-consuming box operations and cumbersome hyper-parameter tuning. Dynamic label assignment methods [15, 47, 70] are proposed to better define the positive and negative samples for model learning. GFL [34, 35] introduces Quality Focal Loss (QFL) and a Distribution-Guided Quality Predictor to increase the consistency between the classification score and the localization quality. It also models the bounding box representation as a probability distribution so that it can capture the localization ambiguity of the box edges. Recently, attributing to the strong ability of the transformer block to encode expressive features, DETR family [4, 6, 30, 42, 45, 68, 77] has become a new trend in the object detection community.

## 2.2. Knowledge Distillation for Object Detection

Knowledge Distillation (KD) is an effective technique to transfer knowledge from a large-scale teacher model to a small-scale student model. It has been widely studied in the classification task [12, 23, 26, 36, 37, 46, 48, 54, 63, 67, 71, 72], but it is still challenging to distill detection models because of the extreme background ratio. The pioneer work [7] proposes the first distillation framework for object detection by simply combining feature imitation and prediction mimicking. Since then, feature imitation has attracted more and more research attention. Typically, some works [13, 25, 33, 61] focus on selecting effective distillation regions for better feature imitation, while other works [19, 31, 75] aim to weight the imitation loss better. There are also methods [5, 65, 66, 69] attempting to design new teacher-student consistency functions, aiming to explore more consistency information or release the strict limit of the MSE loss.

As the earliest distillation strategy proposed in [24], prediction mimicking plays a vital role in classification distillation. Recently, some improved prediction mimicking methods have been proposed to adapt to object detection. For example, Rank Mimicking [31] regards the score rank of the teacher as a kind of knowledge and aims to force the student to rank instances as the teacher. LD [74] proposes to distill the localization distribution of bounding box [35] to transfer localization knowledge. In this paper, we construct a CrossKD pipeline which separates detection and distillation into different heads to alleviate the target conflict problem of prediction mimicking. It's worth noting that HEAD [59] delivers the student features to an independent assistant head to bridge the gap between heterogeneous teacher-student pairs. In contrast, we observe that simply delivering the student feature to the teacher is effective enough to achieve SOTA results. This makes our method quite concise and different from HEAD. Our method is also related to [1, 28, 32, 64], but all of them aim to distill classification models and are not tailored for object detection.
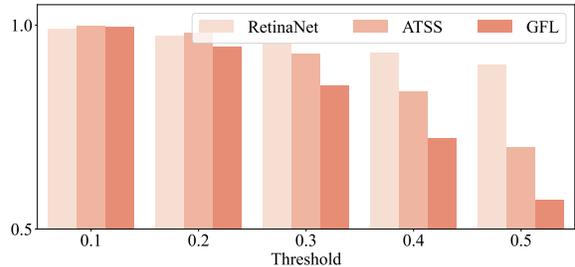


Figure 3. Statistics of the target conflict degree between student (GFL-R50) and teacher (GFL-R101, ATSS-R101, RetinaNet-R101). X-axis is the teacher-student discrepancy threshold for conflict areas. Y-axis represents the ratios of the target conflict areas to the positive areas.

## 3. Methodology

### 3.1. Analysis of the Target Conflict Problem

Target conflict is a common issue confronted in conventional prediction mimicking methods. In contrast to the classification task, which assigns a specific category to every image, the labels in advanced detectors are usually dynamically assigned and not deterministic. Typically, detectors depend on a hand-crafted principle, *i.e.*, assigner, to determine the label in each location. In most cases, detectors cannot reproduce the assigner's labels exactly, which results in a conflict between the teacher-student targets in KD. Furthermore, the inconsistency of the assigners of the student and teacher in real-world scenarios extends the distance between the ground-truth and distillation targets.

To quantitatively measure the degree of target conflict, we statistic the ratios of conflict areas to the positive areas under different teacher-student discrepancy in the COCO *minival* dataset and report the results in Fig. 3. As we can see, even if both the teacher (ATSS [70] and GFL [35]) and student (GFL) have the same label assignment strategy, there are still numerous locations that have a discrepancy larger than 0.5 between the ground-truth and distillation targets, respectively. When we use a teacher with a different assigner (RetinaNet) to distill the student (GFL), the conflict areas increases by a large margin. More experiments in Sec. 4.5 also demonstrate that the target conflict problem severely hinders the performance of prediction mimicking.

Despite the large influence of target conflict, this problem has been neglected for a long time in previous prediction mimicking methods [24, 31]. These methods intend to directly minimize the discrepancy between the teacher-student predictions. Its objective can be described as:

$$\mathcal{L}_{\text{KD}} = \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{R}} \mathcal{S}(r) \mathcal{D}_{\text{pred}}(\boldsymbol{p}^s(r), \boldsymbol{p}^t(r)), \quad (1)$$

where $\boldsymbol{p}^s$ and $\boldsymbol{p}^t$ are the prediction vectors generated by the detection heads of the student and the teacher, respectively.
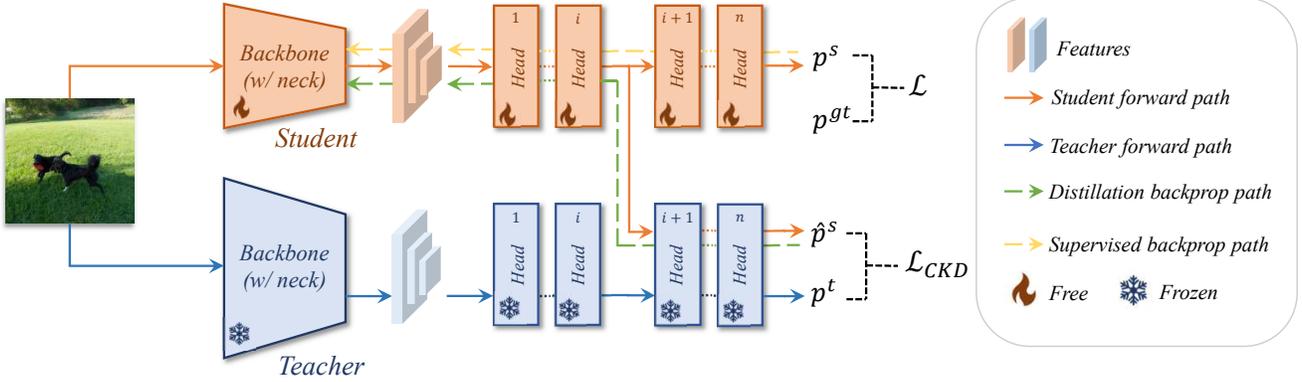
Figure 4. Overall framework of the proposed CrossKD. For a given teacher-student pair, CrossKD first delivers the intermediate features of the student into the teacher layers and generates the cross-head predictions $\hat{\boldsymbol{p}}^s$. Then, distillation losses are calculated between the original teacher's predictions and the cross-head predictions of the student. In back-propagation, the gradients with respect to the detection loss normally pass through the student detection head, while the distillation gradients propagate through the frozen teacher layers.

$\mathcal{D}_{\text{pred}}(\cdot)$ refers to the loss function calculating the discrepancy between $\boldsymbol{p}^s$ and $\boldsymbol{p}^t$, *e.g.*, KL Divergence [24] for classification, L1 Loss [7] and LD [74] for regression. $\mathcal{S}(\cdot)$ is the region selection principle which produces a weight at each position $r$ in the entire image region $\mathcal{R}$.

It's worth noting that $\mathcal{S}(\cdot)$, to a certain extent, can alleviate the target conflict problem by down-weighting the regions with large teacher-student discrepancies. However, the heavily uncertain regions usually accommodate more information benefits for the student than undisputed areas. Ignoring those regions may have a large impact on the effectiveness of prediction mimicking methods. Consequently, to push the envelope of prediction mimicking, it is necessary to handle the target conflict problem gracefully instead of directly down-weighting.

### 3.2. Cross-Head Knowledge Distillation

As described in Sec. 3.1, we observe that directly mimicking the predictions of the teacher confronts the target conflict problem, which hinders prediction mimicking achieving promising performance. To alleviate this problem, we present a novel Cross-head Knowledge Distillation (CrossKD) in this section. The overall framework is illustrated in Fig. 4. Like many previous prediction mimicking methods, our CrossKD performs the distillation process on the predictions. Differently, CrossKD delivers the intermediate features of the student to the teacher's detection head and generates cross-head predictions to conduct distillation.

Given a dense detector, like RetinaNet [39], each detection head usually consists of a sequence of convolutional layers, represented as $\{C_i\}$. For simplicity, we suppose there are totally $n$ convolutional layers in each detection head (e.g., 5 in RetinaNet with 4 hidden layers and 1 prediction layer). We use $\boldsymbol{f}_i, i \in \{1, 2, \cdots, n-1\}$ to denote the feature maps produced by $C_i$ and $\boldsymbol{f}_0$ the input feature

maps of $C_1$. The predictions $\boldsymbol{p}$ are generated by the last convolutional layer $C_n$. Thus, for a given teacher-student pair, the predictions of the teacher and the student can be represented as $\boldsymbol{p}^t$ and $\boldsymbol{p}^s$, respectively.

Besides the original predictions from the teacher and the student, CrossKD additionally delivers the student's intermediate features $\boldsymbol{f}_i^s, i \in \{1, 2, \cdots, n-1\}$ to $C_{i+1}^t$, the $(i+1)$-th convolutional layer of the teacher's detection head, resulting in the cross-head predictions $\hat{\boldsymbol{p}}^s$. Given $\hat{\boldsymbol{p}}^s$, instead of computing the KD loss between $\boldsymbol{p}^s$ and $\boldsymbol{p}^t$, we propose to use the KD loss between the cross-head predictions $\hat{\boldsymbol{p}}^s$ and the original predictions of the teacher $\boldsymbol{p}^t$ as the objective of our CrossKD, which is described as follows:

$$\mathcal{L}_{\text{CrossKD}} = \frac{1}{|\mathcal{S}|} \sum_{r \in \mathcal{R}} \mathcal{S}(r) \mathcal{D}_{\text{pred}}(\hat{\boldsymbol{p}}^s(r), \boldsymbol{p}^t(r)), \tag{2}$$

where $\mathcal{S}(\cdot)$ and $|\mathcal{S}|$ are the region selection principle and the normalization factor. Instead of designing complicated $\mathcal{S}(\cdot)$, we equally conduct distillation between $\hat{\boldsymbol{p}}^s$ and $\boldsymbol{p}^t$ over the entire prediction map. Specifically, $\mathcal{S}(\cdot)$ is a constant function with the value of 1 in our CrossKD. According to the different tasks of each branch (e.g., classification or regression), we perform different types of $\mathcal{D}_{\text{pred}}(\cdot)$ to effectively deliver task-specific knowledge to the student.

By performing CrossKD, the detection loss and the distillation loss are separately applied to different branches. As illustrated in Fig. 4, the gradients of the detection loss pass through the entire head of the student, while the gradients of distillation loss propagate through the frozen teacher layers to the latent features of the student, which heuristically increases the consistency between the teacher and the student. Compared to directly closing the predictions between the teacher-student pair, CrossKD allows part of the student's detection head to be only relative with detection losses, re-

Table 1. Effectiveness of applying CrossKD at different positions. The index $i$ represents the intermediate features used as input in the cross-head branches. 'LD' means the direct application of prediction mimicking on the student's head with LD [74]. The teacher-student pair is GFL with ResNet-50 and ResNet-18 backbones. We can see that $i = 3$ yields the best performance in this experiment.

| $i$ | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|
| - | 35.8 | 53.1 | 38.2 | 18.9 | 38.9 | 47.9 |
| 0 | 38.2 | 55.6 | 41.3 | 20.2 | 41.9 | 50.9 |
| 1 | 38.3 | 55.8 | 41.1 | 20.8 | 42.1 | 49.8 |
| 2 | 38.6 | 56.2 | 41.5 | 20.8 | 42.7 | 50.7 |
| 3 | 38.7 | 56.3 | 41.6 | 21.1 | 42.2 | 51.1 |
| 4 | 38.2 | 55.7 | 41.2 | 20.3 | 41.9 | 50.2 |
| LD | 37.8 | 55.5 | 40.5 | 20.0 | 41.4 | 49.5 |

Table 2. Comparisons between feature imitation and CrossKD. we choose advanced PKD to represent feature imitation and apply PKD to different positions to compare with CrossKD fairly. Here, 'neck' means performing PKD on the FPN neck. 'cls' and 'reg' indicate applying PKD to the classification branch and the regression, respectively. The teacher-student pair is GFL with ResNet-50 and ResNet-18 backbones.

| Methods | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|
| - | 35.8 | 53.1 | 38.2 | 18.9 | 38.9 | 47.9 |
| PKD:neck | 38.0 | 55.0 | 41.2 | 19.6 | 41.5 | 50.2 |
| PKD:cls | 37.5 | 54.9 | 40.5 | 19.5 | 41.1 | 50.5 |
| PKD:reg | 37.2 | 54.0 | 40.2 | 19.0 | 40.9 | 50.0 |
| PKD:cls+reg | 37.3 | 54.3 | 40.0 | 19.2 | 41.1 | 49.8 |
| **CrossKD** | 38.7 | 56.3 | 41.6 | 21.1 | 42.2 | 51.1 |

sulting in a better optimization towards ground-truth targets. Quantitative analysis is presented in our experiment section.

### 3.3. Optimization Objectives

The overall loss for training can be formulated as the weighted sum of the detection loss and the distillation loss, written as:

$$
\begin{aligned}
\mathcal{L} = \ &\mathcal{L}_{\text{cls}}(\boldsymbol{p}_{\text{cls}}^s, \boldsymbol{p}_{\text{cls}}^{gt}) + \mathcal{L}_{\text{reg}}(\boldsymbol{p}_{\text{reg}}^s, \boldsymbol{p}_{\text{reg}}^{gt}) \\
&+ \mathcal{L}_{\text{CrossKD}}^{\text{cls}}(\hat{\boldsymbol{p}}_{\text{cls}}^s, \boldsymbol{p}_{\text{cls}}^t) + \mathcal{L}_{\text{CrossKD}}^{\text{reg}}(\hat{\boldsymbol{p}}_{\text{reg}}^s, \boldsymbol{p}_{\text{reg}}^t),
\end{aligned}
\tag{3}
$$

where $\mathcal{L}_{\text{cls}}$ and $\mathcal{L}_{\text{reg}}$ stand for the detection losses which are calculated between the student predictions $\boldsymbol{p}_{\text{cls}}^s$, $\boldsymbol{p}_{\text{reg}}^s$ and their corresponding ground truth targets $\boldsymbol{p}_{\text{cls}}^{gt}$, $\boldsymbol{p}_{\text{reg}}^{gt}$. The additional CrossKD losses are represented as $\mathcal{L}_{\text{CrossKD}}^{\text{cls}}$ and $\mathcal{L}_{\text{CrossKD}}^{\text{reg}}$, which are performed between the cross-head predictions $\hat{\boldsymbol{p}}_{\text{cls}}^s$, $\hat{\boldsymbol{p}}_{\text{reg}}^s$ and the teacher's predictions $\boldsymbol{p}_{\text{cls}}^t$, $\boldsymbol{p}_{\text{reg}}^t$.

We apply different distance functions $\mathcal{D}_{\text{pred}}$ to transfer task-specific information in different branches. In the classification branch, we regard the classification scores predicted by the teacher as the soft labels and directly use Quality Focal Loss (QFL) proposed in GFL [35] to pull close the teacher-student distance. As for regression, there are mainly two types of regression forms presenting in dense detectors. The first regression form directly regresses the bounding boxes from the anchor boxes (e.g., RetinaNet [39], ATSS [70]) or points (e.g., FCOS [57]). In this case, we directly use GIoU [53] as $\mathcal{D}_{\text{pred}}$. In the other situation, the regression form predicts a vector to represent the distribution of box location (e.g., GFL [35]), which contains richer information than the Dirac distribution of the bounding box representation. To efficiently distill the knowledge of location distribution, we employ KL divergence, like LD [74], to transfer localization knowledge. More details about the loss functions are given in the supplementary materials.

## 4. Experiments

### 4.1. Implement Details

We evaluate the proposed method on the large-scale MS COCO [40] benchmark as done in most previous works. To ensure consistency with the standard practice, we use the *trainval135k* set (115$K$ images) for training and the *minival* set (5$K$ images) for validation. For evaluation, the standard COCO-style measurement, i.e., Average Precision (AP), is used. We also report mAP with IoU thresholds of 0.5 and 0.75, as well as AP for small, medium, and large objects. Our proposed method, CrossKD, is implemented under the MMDetection [9] framework in Python. For a fair comparison, all experiments are developed using 8 Nvidia V100 GPUs with a minibatch of two images per GPU. Unless otherwise stated, all the hyper-parameters follow the default settings of the corresponding student model for both training and testing.

### 4.2. Method Analysis

To investigate the effectiveness of our method, we conduct extensive ablation experiments based on GFL [35]. If not specified, we use GFL with the ResNet-50 backbone [22] as the teacher detector and use the ResNet-18 backbone in the student detector. The accuracy of the teacher and the student are 40.2 AP and 35.8 AP, respectively. All experiments follow the default 1× training schedule (12 epochs).

**Positions to apply CrossKD.** As described in Sec. 3.2, CrossKD delivers the $i$-th intermediate feature of the student to part of the teacher's head. Here, we conduct distillation on both classification and box regression branches. When $i = 0$, CrossKD directly feeds the student's FPN features into the teacher's head. In this case, the entire student's head is only supervised by the detection loss, and no distillation loss is involved. As $i$ gradually increases, more layers of the student's head are jointly affected by the detec-
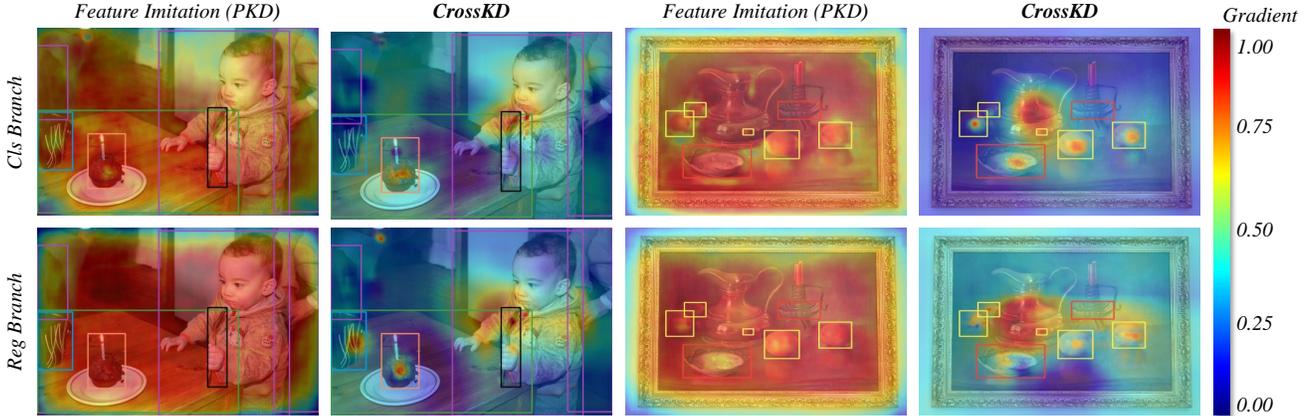
Figure 5. Visualizations of the gradients w.r.t feature imitation and CrossKD. The visualization demonstrates that our CrossKD guided by prediction mimicking can effectively focus on the potentially valuable regions.

Table 3. Effectiveness of CrossKD on different branches. We separately apply CrossKD on the classification (cls) and regression (reg) branches. The teacher-student pair is GFL with ResNet-50 and ResNet-18 backbones.

| cls | reg | LD | | | CrossKD | | |
|---|---|---|---|---|---|---|---|
| | | AP | $AP_{50}$ | $AP_{75}$ | AP | $AP_{50}$ | $AP_{75}$ |
| ✓ | | 37.3 | 55.2 | 40.0 | 37.7 | 55.6 | 40.2 |
| | ✓ | 36.8 | 53.8 | 39.6 | 37.2 | 54.0 | 40.0 |
| ✓ | ✓ | 37.8 | 55.4 | 40.5 | 38.7 | 56.3 | 41.6 |

Table 4. Collective effect of CrossKD and prediction mimicking. The teacher-student pair is GFL with ResNet-50 and ResNet-18 backbones.

| CrossKD | LD | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|---|
| - | - | 35.8 | 53.1 | 38.2 | 18.9 | 38.9 | 47.9 |
| ✓ | | 38.7 | 56.3 | 41.6 | 21.1 | 42.2 | 51.1 |
| | ✓ | 37.8 | 55.5 | 40.5 | 20.0 | 41.4 | 49.5 |
| ✓ | ✓ | 38.1 | 55.6 | 40.9 | 20.4 | 41.6 | 51.1 |

tion loss and the distillation loss. When $i = n$, our method degrades to the original prediction mimicking, where the distillation loss will be directly performed between the two predictions of the teacher-student pair.

In Tab. 1, we report the results of performing CrossKD on different intermediate features. One can see that our CrossKD can improve the distillation performance for all the choices of $i$. Notably, when using the 3-rd intermediate features, CrossKD reaches the best performance of 38.7 AP, which is 0.9 AP higher than the recent state-of-the-art prediction mimicking method LD [74]. This suggests that not all layers in the student's head need to be isolated from the influence of the distillation loss. Therefore, we use $i = 3$ as the default setting in all subsequent experiments.

**CrossKD v.s. Feature Imitation.** We compare CrossKD with the advanced feature imitation method PKD [5]. For a fair comparison, we perform PKD on the same positions as our CrossKD, including FPN features and the third layer of detection heads. The results are reported in Tab. 2. It can be seen that PKD can achieve 38.0 AP when it is applied between the FPN features of the teacher-student pair. On the detection head, PKD even shows a performance drop. In contrast, our CrossKD achieves 38.7 AP, which is 0.7 AP higher than PKD applied on the FPN features.

To further investigate the advantage of CrossKD, we visualize the gradients on the latent features of the detection

head, as shown in Fig. 5. As illustrated, the gradients generated by PKD have a large and wide impact on the entire feature maps, which is inefficient and not targeted. On the contrary, the gradients generated by CrossKD can focus on potential semantic areas with objects of interest.

**CrossKD v.s. Prediction Mimicking.** We begin by separately performing prediction mimicking and CrossKD on the classification and box regression branches. The results are reported in Tab. 3. One can see that replacing prediction mimicking with CrossKD leads to a stable performance gain regardless of classification or regression branches. Specifically, prediction mimicking produces 37.3 AP and 36.8 AP on the classification and regression branches, respectively, while CrossKD yields 37.7 AP and 37.2 AP, representing a consistent improvement over the corresponding results of prediction mimicking. If KD is performed on the two branches, our method can still outperform prediction mimicking by +0.9 AP. Moreover, we further evaluate the collective effect of prediction mimicking and CrossKD, as shown in Tab. 4. Intriguingly, we observe that using both prediction mimicking and CrossKD together yields a final result of 38.1 AP, which is even lower than the result of using CrossKD alone. We believe that this is because the prediction mimicking introduces the target conflict problem again, which makes the student model struggle to learn.

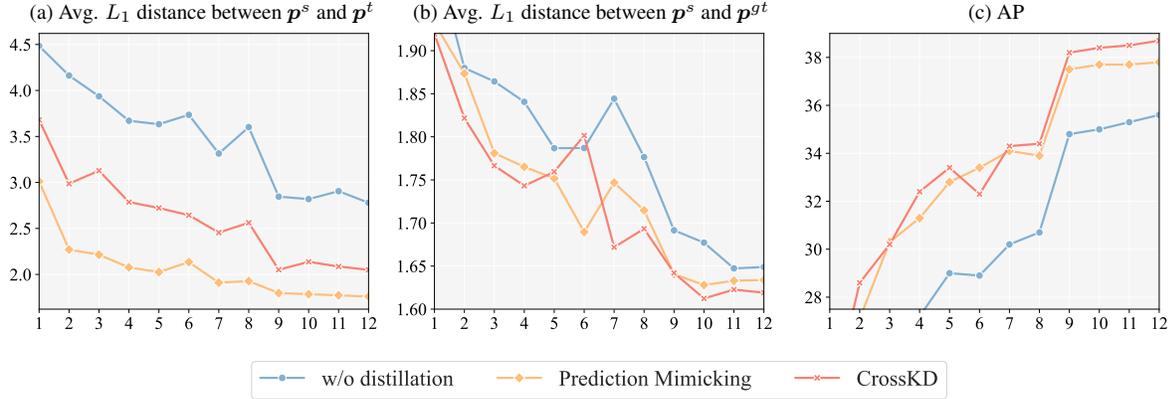In addition, we visualize the statistical variation during

Figure 6. Visualization for the variation of statistics during training. (a) Curves of average $L_1$ distance between student predictions $\boldsymbol{p}^s$ and teacher's $\boldsymbol{p}^t$. (b) Curves of average $L_1$ distance between student predictions $\boldsymbol{p}^s$ and positive ground truth targets $\boldsymbol{p}^{gt}$. (c) Curves of Average Precision (AP). All curves are evaluated on the COCO *minival* set. X-axis refers to the epoch number. Y-axis in (a) and (b) indicate the average $L_1$ distance, while in (c) means the value of AP.

Table 5. Comparison with state-of-the-art detection KD methods on COCO. * denotes results are referenced from LD [74] and PKD [5]. All results are evaluated on the COCO *minival* set.

| Method | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|
| GFL-R101 (T) | 44.9 | 63.1 | 49.0 | 28.0 | 49.1 | 57.2 |
| GFL-R50 (S) | 40.2 | 58.4 | 43.3 | 23.3 | 44.0 | 52.2 |
| FitNets* [54] | 40.7 (0.5↑) | 58.6 | 44.0 | 23.7 | 44.4 | 53.2 |
| Inside GT Box* | 40.7 (0.5↑) | 58.6 | 44.2 | 23.1 | 44.5 | 53.5 |
| Defeat* [19] | 40.8 (0.6↑) | 58.6 | 44.2 | 24.3 | 44.6 | 53.7 |
| Main Region* [74] | 41.1 (0.9↑) | 58.7 | 44.4 | 24.1 | 44.6 | 53.6 |
| Fine-Grained* [62] | 41.1 (0.9↑) | 58.8 | 44.8 | 23.3 | 45.4 | 53.1 |
| FGD [65] | 41.3 (1.1↑) | 58.8 | 44.8 | 24.5 | 45.6 | 53.0 |
| GID* [13] | 41.5 (1.3↑) | 59.6 | 45.2 | 24.3 | 45.7 | 53.6 |
| SKD [14] | 42.3 (2.1↑) | 60.2 | 45.9 | 24.4 | 46.7 | 55.6 |
| LD [74] | 43.0 (2.8↑) | 61.6 | 46.6 | 25.5 | 47.0 | 55.8 |
| PKD* [5] | 43.3 (3.1↑) | 61.3 | 46.9 | 25.2 | 47.9 | 56.2 |
| **CrossKD** | 43.7 (3.5↑) | **62.1** | 47.4 | **26.9** | 48.0 | 56.2 |
| **CrossKD+PKD** | **43.9 (3.7↑)** | 62.0 | **47.7** | 26.4 | **48.5** | **57.0** |

Table 6. CrossKD for detectors with homogeneous backbones. Teacher detectors use ResNet-101 as the backbone, while the students use ResNet-50 as the backbone. All results are evaluated on the COCO *minival* set.

| Student | Methods | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|---|
| RetinaNet [39] | R101 | 38.9 | 58.0 | 41.5 | 21.0 | 32.8 | 52.4 |
| | R50 | 37.4 | 56.7 | 39.6 | 20.0 | 40.7 | 49.7 |
| | **CrossKD** | 39.7 | 58.9 | 42.5 | 22.4 | 43.6 | 52.8 |
| FCOS [57] | R101 | 40.8 | 60.0 | 44.0 | 24.2 | 44.3 | 52.4 |
| | R50 | 38.5 | 57.7 | 41.0 | 21.9 | 42.8 | 48.6 |
| | **CrossKD** | 41.3 | 60.6 | 44.2 | 25.1 | 45.5 | 52.4 |
| ATSS [70] | R101 | 41.5 | 59.9 | 45.2 | 24.2 | 45.9 | 53.3 |
| | R50 | 39.4 | 57.6 | 42.8 | 23.6 | 42.9 | 50.3 |
| | **CrossKD** | 41.8 | 60.1 | 45.4 | 24.9 | 45.9 | 54.2 |

## 4.3. Comparison with SOTA KD Methods

In this section, we evaluate various state-of-the-art object detection KD methods on the GFL [35] framework and fairly compare them with our proposed CrossKD. We use ResNet-101 as the backbone for the teacher detector, which is trained with a $2\times$ schedule and multi-scale augmentation. For the student detector, we adopt the ResNet-50 backbone. We train the student with the $1\times$ schedule. The pretrained checkpoint of the teacher is directly borrowed from the MMDetection[9] model zoo.

We report all results in Tab. 5. As we can see, at the same condition, CrossKD can achieve 43.7 AP without bells and whistles, which improves the accuracy of the student by 3.5 AP, outperforming all other state-of-the-art methods. Notably, CrossKD surpasses the advanced feature imitation method PKD by 0.4 AP and surpasses the advanced prediction mimicking method LD by 0.7 AP, demonstrating the effectiveness of CrossKD. In addition, we also observe that

training to conduct further analysis on CrossKD and prediction mimicking. We first calculate the $L_1$ distances between the student's predictions $\boldsymbol{p}^s$ and the teacher's predictions $\boldsymbol{p}^t$, as well as the ground-truth targets $\boldsymbol{p}^{gt}$ at each epoch. As plotted in Fig. 6(a), the distance $L_1(\boldsymbol{p}^s, \boldsymbol{p}^t)$ can be reduced significantly by our CrossKD, while it is reasonable for the prediction mimicking to achieve the lowest distance as the distillation is directly imposed on $\boldsymbol{p}^s$. However, as the optimization target conflict exists, the prediction mimicking involves a contradictory optimization process, thereby generally yielding a larger distance $L_1(\boldsymbol{p}^s, \boldsymbol{p}^{gt})$ than our CrossKD, as shown in Fig. 6(b). In Fig. 6(c), our method shows a faster training process and achieves the best performance of 37.8 AP.

Table 7. CrossKD for teacher-student pairs with different label assigners. All results are evaluated on the COCO *minival* set.

| Methods | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|
| GFL-R50 (S) | 40.2 | 58.4 | 43.3 | 23.3 | 44.0 | 52.2 |
| ATSS[70]-R101 (T) | 41.5 | 59.9 | 45.2 | 24.2 | 45.9 | 53.3 |
| KD | 39.7 | 57.9 | 42.8 | 21.8 | 44.2 | 51.5 |
| **CrossKD** | 42.1 | 60.5 | 45.7 | 24.5 | 46.3 | 54.5 |
| GFL-R50 (S) | 40.2 | 58.4 | 43.3 | 23.3 | 44.0 | 52.2 |
| Retinanet[39]-R101 (T) | 38.9 | 58.0 | 41.5 | 21.0 | 32.8 | 52.4 |
| KD | 30.3 | 49.2 | 31.2 | 20.0 | 38.1 | 34.4 |
| **CrossKD** | 41.2 | 59.4 | 44.8 | 24.0 | 45.1 | 53.5 |

Table 8. CrossKD for other detector pairs with Heterogeneous Backbones. For convenience, only the backbone lists below, where 'SwinT' refers to RetinaNet with a tiny version of Swin Transformer [43]. All results are evaluated on the COCO *minival* set.

| Methods | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|
| SwinT (T) [43] | 37.3 | 57.5 | 39.9 | 22.7 | 41.0 | 49.6 |
| ResNet-50 (S) | 36.5 | 55.4 | 39.1 | 20.4 | 40.3 | 48.1 |
| PKD | 37.2 | 56.7 | 39.5 | 21.2 | 41.2 | 49.7 |
| **CrossKD** | 38.0 | 58.1 | 40.5 | 23.1 | 41.8 | 49.7 |
| ResNet-50 (T) | 36.5 | 55.4 | 39.1 | 20.4 | 40.3 | 48.1 |
| MobileNetv2 (S) [55] | 30.9 | 48.7 | 32.5 | 16.3 | 33.5 | 41.9 |
| PKD | 33.2 | 51.3 | 35.0 | 16.5 | 36.6 | 46.5 |
| **CrossKD** | 34.1 | 52.7 | 36.5 | 18.8 | 37.1 | 45.4 |

CrossKD is also orthogonal to the feature imitation methods. With the help of PKD, CrossKD achieves the highest results of 43.9 AP, achieving an improvement of 3.7 AP compared to the baseline.

## 4.4. CrossKD on Different Detectors

Besides performing CrossKD on GFL, we select three commonly used detectors, i.e., RetinaNet[39], FCOS [57], and ATSS [70], to investigate the effectiveness of CrossKD. We strictly follow the student settings for training and reference the teacher and student results from the MMDetection model zoo. The results are presented in Tab. 6. As shown in Tab. 6, CrossKD significantly boosts the performance of all three types of detectors. Specifically, RetinaNet, FCOS, and ATSS with our CrossKD achieve 39.7 AP, 41.3 AP, and 41.8 AP, respectively, which are 2.3 AP, 2.8 AP, and 2.4 AP higher than their corresponding baselines. All results after distillation even outperform the original teachers, indicating that CrossKD can work well on different dense detectors.

## 4.5. Distillation under Severe Target Conflict

In this subsection, we perform prediction mimicking and our CrossKD between detectors with different assigners to explore the effectiveness of CrossKD against the target conflict problem. As shown in Tab. 7, the target conflict problem has a large impact on the optimization of the student, leading to an inferior performance. Specifically, prediction mimicking reduces the AP to 30.3 with the teacher as RetinaNet which has a different assigner with GFL. Furthermore, even if the ATSS has the same assigner as GFL, the student's AP is only distilled to 39.7, falling below the performance without KD. In contrast, CrossKD can still significantly improve the student's accuracy even if existing a large discrepancies between the ground-truth and distillation targets. CrossKD boosts the accuracy of GFL-R50 to 42.1 (+1.9 AP) when applying ATSS as the teacher. Even guided by a weak teacher ReitnaNet, CrossKD still improves the performance of GFL-R50 to 41.2 AP, 1.0 AP higher than the baseline. This demonstrates how robust our

CrossKD is when confronting severe target conflict.

## 4.6. Distillation between Heterogeneous Backbones

In this subsection, we explore the ability of our CrossKD for distilling the heterogeneous students. We perform knowledge distillation on RetinaNet [39] with different backbone networks and compare our method with the recent state-of-the-art method PKD [5]. Specifically, we choose two typical heterogeneous backbones, i.e., the transformer backbone Swin-T [43] and the lightweight backbone MobileNetv2 [55]. All the detectors are trained for 12 epochs with a single-scale strategy. The results are presented in Tab. 8. We can see when distilling knowledge from Swin-T, CrossKD reaches 38.0 AP (+1.5 AP), outperforming PKD by 0.8 AP. CrossKD also improves the results of RetinaNet with the MoblieNetv2 backbone to 34.1 AP, which is 3.2 AP higher than the baseline and surpasses PKD by 0.9 AP.

## 5. Conclusions and Discussions

In this paper, we introduce CrossKD, a novel KD method designed to enhance the performance of dense object detectors. CrossKD transfers the intermediate features from the student's head to that of the teacher to produce the cross-head predictions for distillation, an efficient way to alleviate the conflict between the supervised and distillation targets. Our results have shown that CrossKD can improve the distillation efficiency and achieve state-of-the-art performance. In the future, we will further extend our method to other relevant fields, *e.g.* 3D object detection.

# References

[1] Haoli Bai, Jiaxiang Wu, Irwin King, and Michael Lyu. Few shot network compression via cross distillation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.

[2] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934*, 2020.

[3] Qinglong Cao, Zhengqin Xu, Yuantian Chen, Chao Ma, and Xiaokang Yang. Domain-controlled prompt learning. *arXiv preprint arXiv:2310.07730*, 2023.

[4] Qinglong Cao, Zhengqin Xu, Yuntian Chen, Chao Ma, and Xiaokang Yang. Domain prompt learning with quaternion networks. *arXiv preprint arXiv:2312.08878*, 2023.

[5] Weihan Cao, Yifan Zhang, Jianfei Gao, Anda Cheng, Ke Cheng, and Jian Cheng. Pkd: General distillation framework for object detectors via pearson correlation coefficient. In *Advances in Neural Information Processing Systems*, 2022.

[6] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 213–229. Springer, 2020.

[7] Guobin Chen, Wongun Choi, Xiang Yu, Tony Han, and Manmohan Chandraker. Learning efficient object detection models with knowledge distillation. *Advances in neural information processing systems*, 30, 2017.

[8] Kai Chen, Jiangmiao Pang, Jiaqi Wang, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. Hybrid task cascade for instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[9] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. MMDetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019.

[10] Shaoyu Chen, Tianheng Cheng, Jiemin Fang, Qian Zhang, Yuan Li, Wenyu Liu, and Xinggang Wang. Tinydet: accurately detecting small objects within 1 gflops. *Science China Information Sciences*, 66(1):119102, 2023.

[11] Yuming Chen, Xinbin Yuan, Ruiqi Wu, Jiabao Wang, Qibin Hou, and Ming-Ming Cheng. Yolo-ms: Rethinking multi-scale representation learning for real-time object detection, 2023.

[12] Jang Hyun Cho and Bharath Hariharan. On the efficacy of knowledge distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.

[13] Xing Dai, Zeren Jiang, Zhao Wu, Yiping Bao, Zhicheng Wang, Si Liu, and Erjin Zhou. General instance distillation for object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7842–7851, June 2021.

[14] Philip De Rijk, Lukas Schneider, Marius Cordts, and Dariu Gavrila. Structural knowledge distillation for object detection. In *Advances in Neural Information Processing Systems*, 2022.

[15] Chengjian Feng, Yujie Zhong, Yu Gao, Matthew R. Scott, and Weilin Huang. Tood: Task-aligned one-stage object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3510–3519, October 2021.

[16] Zheng Ge, Songtao Liu, Feng Wang, Zeming Li, and Jian Sun. Yolox: Exceeding yolo series in 2021. *arXiv preprint arXiv:2107.08430*, 2021.

[17] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, December 2015.

[18] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014.

[19] Jianyuan Guo, Kai Han, Yunhe Wang, Han Wu, Xinghao Chen, Chunjing Xu, and Chang Xu. Distilling object detectors via decoupled features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2154–2164, June 2021.

[20] Jianyuan Guo, Kai Han, Han Wu, Chao Zhang, Xinghao Chen, Chunjing Xu, Chang Xu, and Yunhe Wang. Positive-unlabeled data purification in the wild for object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2653–2662, 2021.

[21] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.

[22] 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 (CVPR)*, June 2016.

[23] 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.

[24] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network (2015). *arXiv preprint arXiv:1503.02531*, 2, 2015.

[25] Zihao Jia, Shengkun Sun, Guangcan Liu, and Bo Liu. Mssd: multi-scale self-distillation for object detection. *Visual Intelligence*, 2(1):8, 2024.

[26] Jangho Kim, SeongUk Park, and Nojun Kwak. Paraphrasing complex network: Network compression via factor transfer. *Advances in neural information processing systems*, 31, 2018.

[27] Tao Kong, Fuchun Sun, Huaping Liu, Yuning Jiang, Lei Li, and Jianbo Shi. Foveabox: Beyound anchor-based object detection. *IEEE Transactions on Image Processing*, 29:7389–7398, 2020.

[28] Animesh Koratana, Daniel Kang, Peter Bailis, and Matei Zaharia. LIT: Learned intermediate representation training for model compression. In *International Conference on Learn-*

*ing Representations*, 2019.

[29] Yuqing Lan, Yao Duan, Chenyi Liu, Chenyang Zhu, Yueshan Xiong, Hui Huang, and Kai Xu. Arm3d: Attention-based relation module for indoor 3d object detection. *Computational Visual Media*, 8(3):395–414, 2022.

[30] Feng Li, Hao Zhang, Shilong Liu, Jian Guo, Lionel M. Ni, and Lei Zhang. Dn-detr: Accelerate detr training by introducing query denoising. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13619–13627, June 2022.

[31] Gang Li, Xiang Li, Yujie Wang, Shanshan Zhang, Yichao Wu, and Ding Liang. Knowledge distillation for object detection via rank mimicking and prediction-guided feature imitation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1306–1313, 2022.

[32] Guilin Li, Junlei Zhang, Yunhe Wang, Chuanjian Liu, Matthias Tan, Yunfeng Lin, Wei Zhang, Jiashi Feng, and Tong Zhang. Residual distillation: Towards portable deep neural networks without shortcuts. *Advances in Neural Information Processing Systems*, 33:8935–8946, 2020.

[33] Quanquan Li, Shengying Jin, and Junjie Yan. Mimicking very efficient network for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.

[34] Xiang Li, Wenhai Wang, Xiaolin Hu, Jun Li, Jinhui Tang, and Jian Yang. Generalized focal loss v2: Learning reliable localization quality estimation for dense object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11632–11641, June 2021.

[35] Xiang Li, Wenhai Wang, Lijun Wu, Shuo Chen, Xiaolin Hu, Jun Li, Jinhui Tang, and Jian Yang. Generalized focal loss: Learning qualified and distributed bounding boxes for dense object detection. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 21002–21012. Curran Associates, Inc., 2020.

[36] Zheng Li, Xiang Li, Lingfeng Yang, Borui Zhao, Renjie Song, Lei Luo, Jun Li, and Jian Yang. Curriculum temperature for knowledge distillation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, number 2, pages 1504–1512, 2023.

[37] Zheng Li, Jingwen Ye, Mingli Song, Ying Huang, and Zhigeng Pan. Online knowledge distillation for efficient pose estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11740–11750, 2021.

[38] Tsung-Yi Lin, Piotr Dollar, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.

[39] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollar. Focal loss for dense object detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.

[40] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014.

[41] Dongyang Liu, Meina Kan, Shiguang Shan, and CHEN Xilin. Function-consistent feature distillation. In *International Conference on Learning Representations*, 2019.

[42] Shilong Liu, Feng Li, Hao Zhang, Xiao Yang, Xianbiao Qi, Hang Su, Jun Zhu, and Lei Zhang. DAB-DETR: Dynamic anchor boxes are better queries for DETR. In *International Conference on Learning Representations*, 2022.

[43] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10012–10022, October 2021.

[44] Chengqi Lyu, Wenwei Zhang, Haian Huang, Yue Zhou, Yudong Wang, Yanyi Liu, Shilong Zhang, and Kai Chen. Rtmdet: An empirical study of designing real-time object detectors. *arXiv preprint arXiv:2212.07784*, 2022.

[45] Depu Meng, Xiaokang Chen, Zejia Fan, Gang Zeng, Houqiang Li, Yuhui Yuan, Lei Sun, and Jingdong Wang. Conditional detr for fast training convergence. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3651–3660, October 2021.

[46] Seyed Iman Mirzadeh, Mehrdad Farajtabar, Ang Li, Nir Levine, Akihiro Matsukawa, and Hassan Ghasemzadeh. Improved knowledge distillation via teacher assistant. In *Proceedings of the AAAI conference on artificial intelligence*, pages 5191–5198, 2020.

[47] Chuong H. Nguyen, Thuy C. Nguyen, Tuan N. Tang, and Nam L.H. Phan. Improving object detection by label assignment distillation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 1005–1014, January 2022.

[48] 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 (CVPR)*, June 2019.

[49] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.

[50] Joseph Redmon and Ali Farhadi. Yolo9000: Better, faster, stronger. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.

[51] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.

[52] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015.

[53] Hamid Rezatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[54] 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.

[55] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[56] Lin Song, Jin-Fu Yang, Qing-Zhen Shang, and Ming-Ai Li. Dense face network: A dense face detector based on global context and visual attention mechanism. *Machine Intelligence Research*, 19(3):247–256, 2022.

[57] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.

[58] Jiaqi Wang, Kai Chen, Shuo Yang, Chen Change Loy, and Dahua Lin. Region proposal by guided anchoring. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[59] Luting Wang, Xiaojie Li, Yue Liao, Zeren Jiang, Jianlong Wu, Fei Wang, Chen Qian, and Si Liu. Head: Hetero-assists distillation for heterogeneous object detectors. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part IX*, pages 314–331. Springer, 2022.

[60] Luequan Wang, Hongbin Xu, and Wenxiong Kang. Mvcontrast: Unsupervised pretraining for multi-view 3d object recognition. *Machine Intelligence Research*, 20(6):872–883, 2023.

[61] Tao Wang, Li Yuan, Xiaopeng Zhang, and Jiashi Feng. Distilling object detectors with fine-grained feature imitation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[62] Tao Wang, Li Yuan, Xiaopeng Zhang, and Jiashi Feng. Distilling object detectors with fine-grained feature imitation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[63] Chenglin Yang, Lingxi Xie, Chi Su, and Alan L. Yuille. Snapshot distillation: Teacher-student optimization in one generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[64] Jing Yang, Brais Martinez, Adrian Bulat, Georgios Tzimiropoulos, et al. Knowledge distillation via softmax regression representation learning. In *International Conference on Learning Representations*, 2021.

[65] 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 (CVPR)*, pages 4643–4652, June 2022.

[66] Lewei Yao, Renjie Pi, Hang Xu, Wei Zhang, Zhenguo Li, and Tong Zhang. G-detkd: Towards general distillation framework for object detectors via contrastive and semantic-guided feature imitation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3591–3600, October 2021.

[67] Junho Yim, Donggyu Joo, Jihoon Bae, and Junmo Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.

[68] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel Ni, and Heung-Yeung Shum. DINO: DETR with improved denoising anchor boxes for end-to-end object detection. In *The Eleventh International Conference on Learning Representations*, 2023.

[69] Linfeng Zhang and Kaisheng Ma. Improve object detection with feature-based knowledge distillation: Towards accurate and efficient detectors. In *International Conference on Learning Representations*, 2021.

[70] Shifeng Zhang, Cheng Chi, Yongqiang Yao, Zhen Lei, and Stan Z. Li. Bridging the gap between anchor-based and anchor-free detection via adaptive training sample selection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[71] Ying Zhang, Tao Xiang, Timothy M. Hospedales, and Huchuan Lu. Deep mutual learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[72] Borui Zhao, Quan Cui, Renjie Song, Yiyu Qiu, and Jiajun Liang. Decoupled knowledge distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11953–11962, June 2022.

[73] Wangbo Zhao, Jiasheng Tang, Yizeng Han, Yibing Song, Kai Wang, Gao Huang, Fan Wang, and Yang You. Dynamic tuning towards parameter and inference efficiency for vit adaptation, 2024.

[74] Zhaohui Zheng, Rongguang Ye, Ping Wang, Dongwei Ren, Wangmeng Zuo, Qibin Hou, and Ming-Ming Cheng. Localization distillation for dense object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9407–9416, June 2022.

[75] Du Zhixing, Rui Zhang, Ming Chang, Shaoli Liu, Tianshi Chen, Yunji Chen, et al. Distilling object detectors with feature richness. *Advances in Neural Information Processing Systems*, 34:5213–5224, 2021.

[76] Chenchen Zhu, Yihui He, and Marios Savvides. Feature selective anchor-free module for single-shot object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[77] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable {detr}: Deformable transformers for end-to-end object detection. In *International Conference on Learning Representations*, 2021.

# Appendix

## 6. Details of Distillation Losses

According to the task of detection heads, *i.e.*, classification, and regression, we apply different distance functions $\mathcal{D}_{\text{pred}}$ to transfer task-specific information in different branches. In this section, we introduce the details of distance functions $\mathcal{D}_{\text{pred}}$ applied in CrossKD.

**Regression Branch.** There are mainly two types of regression branches that existed in dense detectors. The first regression branch directly regresses the bounding boxes from the anchor boxes (*e.g.*, RetinaNet [39], ATSS [70]) or points (*e.g.*, FCOS [57]). In this case, we directly use GIoU [53] as $\mathcal{D}_{\text{pred}}$, which is defined as:

$$\mathcal{D}_{\text{pred}}(\mathcal{B}, \mathcal{B}') = \frac{|\mathcal{B} \cap \mathcal{B}'|}{|\mathcal{B} \cup \mathcal{B}'|} - \frac{|\mathcal{C} \setminus (\mathcal{B} \cup \mathcal{B}')|}{|\mathcal{C}|}, \qquad (4)$$

where $\mathcal{B}$ and $\mathcal{B}'$ represent the predicted and ground-truth bounding boxes and $\mathcal{C}$ is the smallest enclosing convex object for $\mathcal{B}$ and $\mathcal{B}'$.

In the other situation, the regression branch predicts a vector to represent the distribution of box location (*e.g.*, GFL [35]), which contains richer information than the Dirac distribution of the bounding box representation. To efficiently distill the knowledge of location distribution, we employ the same $\mathcal{D}_{\text{pred}}$ like LD [74], which is defined as:

$$\mathcal{D}_{\text{pred}}(\boldsymbol{p}, \boldsymbol{p}') = \text{KL}(s(\boldsymbol{p}/\tau), s(\boldsymbol{p}'/\tau)), \qquad (5)$$

where KL means KL divergence, $s(\cdot)$ indicates the Softmax function, and $\tau$ is a factor to smooth the distribution.

**Classification Branch.** Distillation in the classification branch severely suffers from the imbalance of the foreground and background instances problem. To avoid training crash, previous prediction mimicking methods usually design complicated region selection principle to choose effective areas. In contrast, without selecting effective regions, we regard the classification scores predicted by the teacher as the soft labels and directly use Quality Focal Loss (QFL) proposed in GFL [35] to pull close the teacher-student distance. We define $\mathcal{D}_{\text{pred}}$ in the classification branch as:

$$\mathcal{D}_{\text{pred}}(\boldsymbol{p}, \boldsymbol{p}') = (|\sigma(\boldsymbol{p}) - \sigma(\boldsymbol{p}')|)^{\gamma} \cdot \text{BCE}(\sigma(\boldsymbol{p}), \sigma(\boldsymbol{p}')), \quad (6)$$

where $\sigma$ denotes the sigmoid function and BCE indicates binary cross entropy. $(|\sigma(\boldsymbol{p}) - \sigma(\boldsymbol{p}')|)^{\gamma}$ serves as a modulating factor added to the cross entropy function, with a tunable focusing parameter $\gamma \geq 0$. Here, $\gamma$ is set as 1 in all experiments, which we find is the optimum.

We also compare the performance of QFL with the widely used BCE loss. As shown in Tab. 9, The BCE loss

| Loss | Region | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|------|--------|------|------|------|------|------|------|
| - | - | 35.8 | 53.1 | 38.2 | 18.9 | 38.9 | 47.9 |
| BCE | P | 36.3 | 53.8 | 39.1 | 19.1 | 39.6 | 48.3 |
| BCE | N | 36.2 | 53.5 | 38.9 | 19.3 | 40.0 | 48.2 |
| BCE | P+N | 36.9 | 54.3 | 39.5 | 20.0 | 40.7 | 48.4 |
| QFL | P+N | 38.7 | 56.3 | 41.6 | 21.1 | 42.2 | 51.5 |

Table 9. Effectiveness of different distillation losses in classification branch. 'BCE' and 'QFL' means the binary cross entropy loss and quality focal loss, respectively. 'P' and 'N' refer to the positive and negative regions. The teacher-student pair is GFL with ResNet-50 and ResNet-18 backbones.

| Student | CrossKD | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---------|---------|------|------|------|------|------|------|
| ResNet-18 | | 35.8 | 53.1 | 38.2 | 18.9 | 38.9 | 47.9 |
| | ✓ | 39.2 | 57.0 | 42.2 | 22.7 | 43.0 | 51.3 |
| ResNet-34 | | 38.9 | 56.6 | 42.2 | 21.5 | 42.8 | 51.4 |
| | ✓ | 42.4 | 60.4 | 45.8 | 24.4 | 46.8 | 55.6 |
| ResNet-50 | | 40.2 | 58.4 | 43.3 | 23.3 | 44.0 | 52.2 |
| | ✓ | 43.7 | 62.1 | 47.4 | 26.9 | 48.0 | 56.2 |

Table 10. Quantitative results of CrossKD for lightweight detectors. Standard $1\times$ schedule is applied in all experiments. The teacher detector is GFL with ResNet-101 backbones.

can receive 36.3 and 36.2 AP when separately applied on the positive and negative regions. When we perform distillation on both positive and negative regions, BCE loss can only achieve 36.9 AP, far below 38.7 AP of QFL, which demonstrates the effectiveness of the current distillation losses.

## 7. The Generalization Ability of CrossKD

CrossKD is adaptable for any detector distillation since the target conflict is a common problem of object detection distillation due to imperfect teacher predictions. To demonstrate the generalization, we apply CrossKD on detectors with various types of backbones and structures.

The results of our CrossKD on a series of lightweight students distilled with GFL with ResNet-18, ResNet-34, and ResNet-50 backbones are presented in Tab. 10. We apply ResNet-101 as the backbone for the teacher detector. As shown in Tab. 10, our method can effectively enhance the performance of all given lightweight detectors. Specifically, CrossKD achieves stable improvements for the students with ResNet-18, ResNet-34, and ResNet-50 backbones, which reach 39.2 AP, 42.4 AP, and 43.7 AP.

Furthermore, we adapt CrossKD to typical Faster R-CNN (two-stage) and Deformable DETR (DETR-like) detectors and report their performance in Tab. 11. In Faster R-CNN, we deliver the student region features to the R-
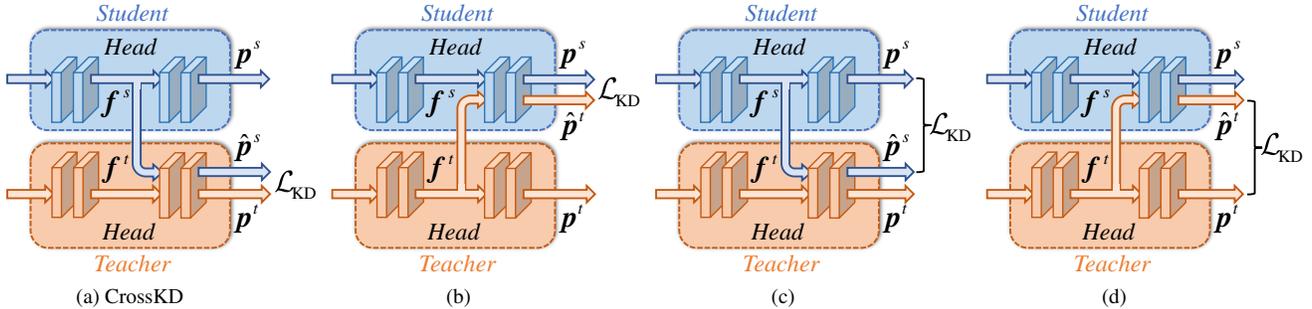
Figure 7. Different cross-head strategies. (a) is the original strategy used in CrossKD. (b) delivers the intermediate features of the teacher to the student head and conducts KD between the cross-head predictions of the teacher and the student's predictions. (c) does the same cross-head strategy as (a) but performs KD between the student's original predictions and cross-head predictions. (d) does the same cross-head strategy as (b) but performs KD between the teacher's original predictions and the cross-head predictions.

CNN head of the teacher to generate cross-head predictions to accept the teacher's supervision. In Deformable DETR, the cross-head predictions are created by passing the encoder features of the student into each stage of the teacher decoder. As shown in Tab. 11, without finely tuned hyper-parameters, CrossKD boosts the accuracy of ResNet-18 based Faster R-CNN and Deformable DETR to 35.5 (2.0 ↑) and 45.8 (1.7↑) AP, which demonstrates the generalization ability of CrossKD.

| Method | Schedule | AP | $AP_{50}$ | $AP_{75}$ |
|---|---|---|---|---|
| Faster R-CNN R18 (S) | 12e | 33.5 | 53.7 | 35.9 |
| Faster R-CNN R50 (T) | 12e | 37.4 | 58.1 | 40.4 |
| CrossKD | 12e | **35.5** (2.0↑) | 55.8 | 38.0 |
| Deform. DETR R18 (S) | 50e | 44.1 | 62.8 | 47.9 |
| Deform. DETR R50 (T) | 50e | 47.0 | 66.1 | 50.9 |
| CrossKD | 50e | **45.8** (1.7↑) | 63.8 | 49.9 |

Table 11. CrossKD for Faster R-CNN and Deformable DETR.

## 8. More Ablations

| Strategy | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|
| - | 35.8 | 53.1 | 38.2 | 18.9 | 38.9 | 47.9 |
| (a) | 38.7 | 56.3 | 41.6 | 21.1 | 42.2 | 51.5 |
| (b) | 35.4 | 52.5 | 37.8 | 18.6 | 38.4 | 47.1 |
| (c) | 34.5 | 51.9 | 36.7 | 17.8 | 37.6 | 45.1 |
| (d) | 32.5 | 48.8 | 35.0 | 16.6 | 35.0 | 42.8 |

Table 12. Comparisons of different cross-head strategies. The strategies (a), (b), (c), (d) have shown in Fig. 7, where (a) is the current strategy used in CrossKD. The teacher-student pair is GFL with ResNet-50 and ResNet-18 backbones.

In this section, we experiment different cross-head

strategies to demonstrate the effectiveness of our CrossKD, which are illustrated in Fig. 7. As presented in Tab. 12, strategy (b), which differently reuses the student's detection head, achieved only 35.4 AP, significantly lower than the 38.7 AP obtained by CrossKD. We hypothesize that this difference in performance may be attributed to the suboptimal optimization of the student's blocks in this approach. Fig. 7(c) and Fig. 7(d) minimize the distances between the original predictions and the cross-head predictions. However, these strategies have limited impact on the student's backbones, resulting in 34.5 AP and 32.5 AP for Fig. 7(c) and Fig. 7(d), respectively.

Moreover, Fig. 7(b), (c), and (d) all perform distillation losses and detection losses at the student's detection heads, so the target conflict problem still exists. In contrast, CrossKD separates the distillation losses onto the teacher's branch and hence avoids the target conflict problem. As a result, CrossKD receives the highest AP of 38.7 among all cross-head strategies.

## 9. Relation to Previous Works

In this section, we describe the differences of our method and some related works which are originally designed for the classification task [1, 28, 32, 41, 64]. Here, we compare CrossKD with these works from the aspects of motivation and structure to emphasize the differences.

**Motivation.** Previous works all concentrate on the classification task. For instance, Bai *et al.* [1] aims to alleviate overfitting in few-shot task. Li *et al.* [32] focuses on using a residual network to help a non-residual network overcome gradient vanishing. Some works [41, 64] target on the general KD scenario in classification. These methods all attempts to solve specific problems in classification and are not specially designed for distilling object detectors.

In contrast, CrossKD, which is specially designed for the object detection task, focuses on the target conflict problem in object detection. To our knowledge, this is the first

work to discuss the target conflict problem in distilling object detectors. As presented in Sec. 1 of the main paper, the teacher detector usually predicts inaccurate results, which conflict with the ground-truth targets. The traditional KD methods supervise the student detector with those two controversial labels at the same place, resulting in low distillation efficiency. To alleviate this problem, we propose to deliver the intermediate features of the student to the part of the teacher's detection head and generate new cross-head predictions to accept the distillation losses.

However, without the detection-specific design, those methods can not achieve a promising performance.

**Structure.** Previous works tend to design a complicated manner to utilize the teacher-student latent features. Typically, Li *et al*. [32] forwards every stage features of the student into the teacher's blocks. Liu *et al*. [41] alternately delivers the intermediate features from the student to the teacher or from the teacher to the student. These strategies significantly increase the computational complexity in training phase.

Instead of applying a complicated design, CrossKD is relatively simple, which only passes the student's latent features through part of the teacher's detection head. Despite its simplicity, extensive experiments demonstrate its effectiveness in object detection KD.

## 10. Result Visualization

We visualize the detection results of the teacher, the student, and our CrossKD in Fig. 8. As the visualization shows, CrossKD usually receives even better results than the teacher detector, which demonstrates that CrossKD can relieve the influence of the teacher's inaccurate predictions and achieve a better optimization towards ground-truths.
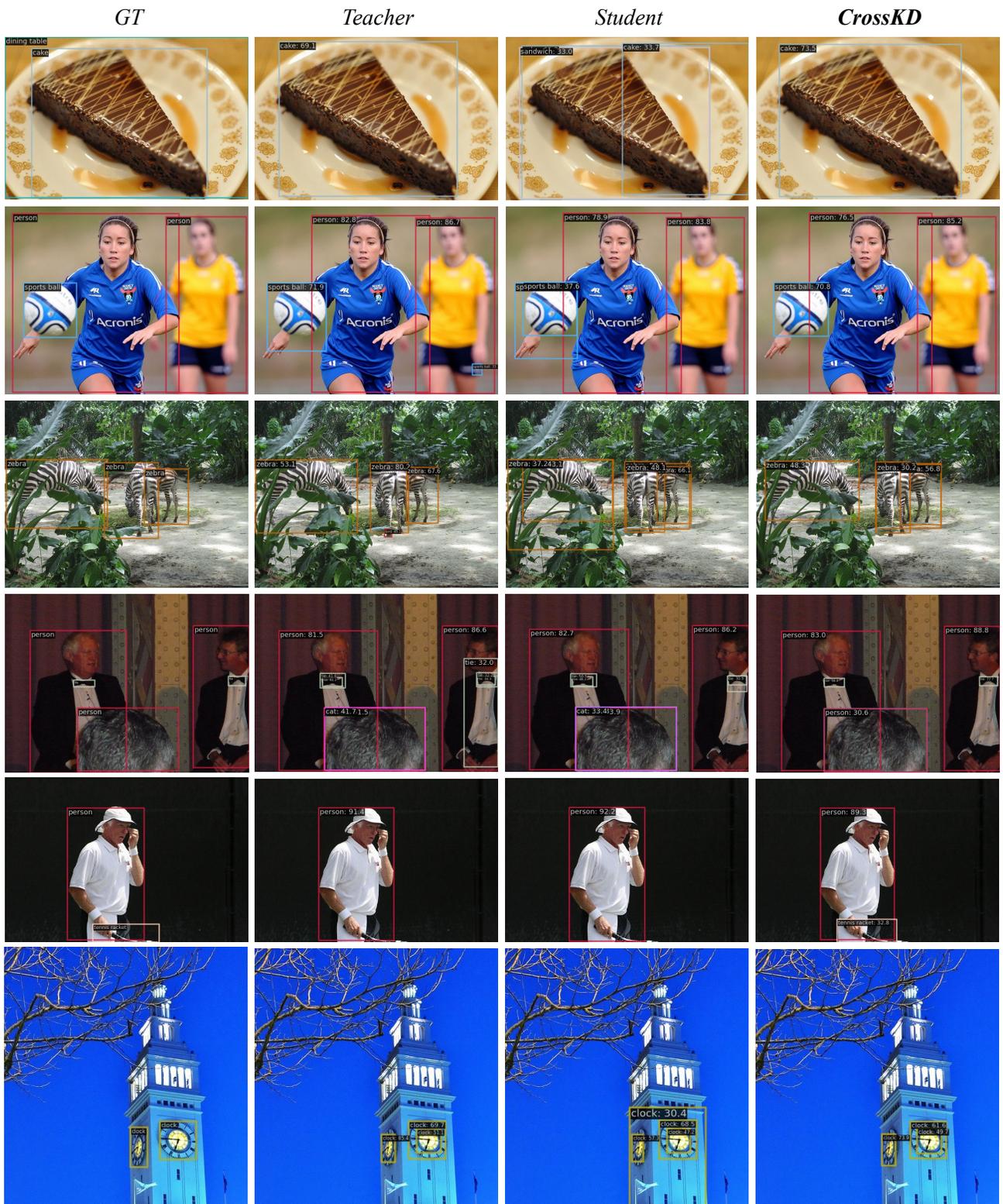
Figure 8. Visualization of detection results of CrossKD. The teacher is GFL-R50 with 40.2 AP and student is GFL-R18 with 35.8 AP.