

Distilling Object Detectors with Global Knowledge

Sanli Tang^{1*}, Zhongyu Zhang^{1*}, Zhanzhan Cheng^{1†}, Jing Lu¹, Yunlu Xu¹, Yi Niu¹, and Fan He²

¹ Hikvision Research Institute, Hanzhou, China

² Institute of Image Processing and Pattern Recognition, Shanghai Jiao Tong University, Shanghai, China

{tangsanli,zhangzhongyu,chengzhanzhan,lujing6,xyunlu,niuyi}@hikvision.com
hf-inspire@sjtu.edu.cn

Abstract. Knowledge distillation learns a lightweight student model that mimics a cumbersome teacher. Existing methods regard the knowledge as the feature of each instance or their relations, which is the instance-level knowledge only from the teacher model, i.e., the *local* knowledge. However, the empirical studies show that the *local* knowledge is much noisy in object detection tasks, especially on the blurred, occluded, or small instances. Thus, a more intrinsic approach is to measure the representations of instances w.r.t. a group of *common* basis vectors in the two feature spaces of the teacher and the student detectors, i.e., *global* knowledge. Then, the distilling algorithm can be applied as space alignment. To this end, a novel prototype generation module (PGM) is proposed to find the *common* basis vectors, dubbed *prototypes*, in the two feature spaces. Then, a robust distilling module (RDM) is applied to construct the global knowledge based on the prototypes and filtrate noisy local knowledge by measuring the discrepancy of the representations in two feature spaces. Experiments with Faster-RCNN and RetinaNet on PASCAL and COCO datasets show that our method achieves the best performance for distilling object detectors with various backbones, which even surpasses the performance of the teacher model. We also show that the existing methods can be easily combined with global knowledge and obtain further improvement. Code is available: <https://github.com/hikvision-research/DAVAR-Lab-ML>.

Keywords: Object Detection, Knowledge Distillation

1 Introduction

Object detectors can be enhanced by applying larger networks [13,21], which, however, will increase the storage and computational cost. A promising solution for finding the sweet spot between efficiency and performance is knowledge distillation (KD) [1,16], which learns a lightweight student that mimics the behaviors of a cumbersome teacher.

* Authors contributed equally. † Corresponding authors.

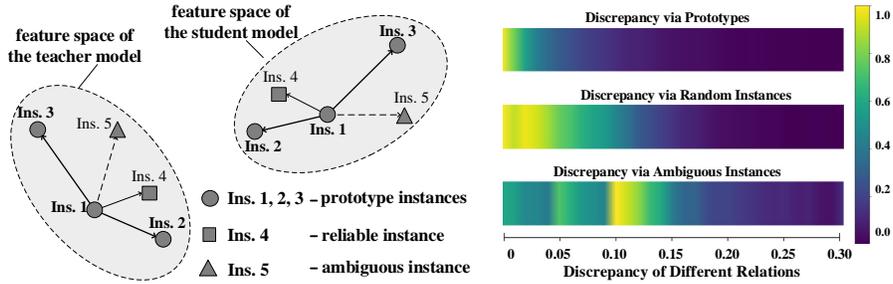


Fig. 1. **Left:** the prototypes are representative and play roles as a common group of basis vectors in TS -space. Although the absolute location of *Ins. 4* is different in TS -space, its representations, e.g., the relations, w.r.t. prototypes are similar while *Ins. 5* shows the representation of much dissimilar. **Right:** on COCO dataset [28] with Faster-RCNN detector [36], we show the discrepancy of relations between instances and three types of basis in TS -space. 10 instances are selected for each class as the bases and others are used for measuring the discrepancy of relations in TS -space. The relations with *prototypes* show much smaller discrepancy than others.

The knowledge can be known to be formed in three categories [10]: feature-based knowledge [37,45,15,43,46], response-based knowledge [16,24,34], and relation-based knowledge [34,30,42,25,4]. Such knowledge can be treated as the *local* knowledge, since only the instance-level knowledge from a single feature space, e.g., the teacher’s, is considered. Based on these knowledge, existing methods design their distilling algorithms for object detection tasks based on some prior senses, e.g., the foreground regions [43], the decoupled background regions [11], the attention guided regions [46,22], or the discrepancy regions [4,23]. However, we find that the local knowledge is of much discrepancy between the teacher and the student in object detection tasks, especially on the ambiguous instances which are blur, truncated, or small. This is because features of ambiguous instances are susceptible to the small disturbance in feature spaces of the teacher and the student. Thus, the distilling process will suffer from the noisy local knowledge, e.g., the false positives and the localization errors, and lead to sub-optimal.

The main concerns on relieving the effect of noisy local knowledge are two folds: constructing reliable global knowledge and applying robust distilling algorithms. By viewing knowledge as the representation of feature space, a more intrinsic approach is to find a group of common basis vectors in both the feature spaces of the teacher and the student detectors. In this way, the *global* knowledge can be formed by representing the instances w.r.t. these basis vectors. Then, a more robust distilling algorithm can be designed by measuring the discrepancy of the representations in the two feature spaces. Hereafter, we name the two feature spaces of the teacher and the student detector as TS -space and the common basis vectors of the TS -space as *prototypes*.

In Fig. 1 (left), we illustrate that: (1) the representations of normal instances w.r.t. the prototypes are of the little discrepancy between two feature spaces, e.g.,

the *Ins.4*; (2) the discrepancy of the ambiguous instances is much larger than others, e.g., the *Ins. 5*. In Fig. 1 (right), we show the statistic analysis of the discrepancy of the instance representations in *TS-space* on the COCO dataset. Notice that each instance is represented by a pair of features in the *TS-space*. Thus, we first measure the cosine similarity between the bases and each of the other instances in the *TS-space*, and then calculate the discrepancy by l_1 distance as shown by the abscissa. In Fig. 1 (right), the discrepancy of relations between prototypes and other instances is much smaller than other bases, which shows a more promising representation of the knowledge in *TS-space*.

Based on the above considerations, we first propose a prototype generation module (PGM) to find a group of common basis vectors as the prototypes in *TS-space*. It selects the prototypes according to minimizing the reconstruction errors of the instances in the two feature spaces, which is inspired by the dictionary learning [41,20,33]. Then, a robust distillation module (RDM) is designed for robust knowledge construction and transfer. Based on the prototypes, the global knowledge is formed by representing the instances under the prototypes, which shows a smaller gap between the two spaces as in Fig. 1 (right). The discrepancy of the representations in *TS-space* can also be regarded as an ensemble of the two models to mitigate noisy local knowledge transferring when distilling. Experiments are carried out with both single-stage (RetinaNet [27]) and two-stage detectors (Faster R-CNN [36]) on Pascal VOC [7] and COCO [28] benchmarks. Extensive experimental results show that the proposed method can effectively improve the performance of knowledge distillation, which achieves new remarkable performance. We also show the existing methods can be further improved by the prototypes with global and local knowledge.

2 Related Works

2.1 Object Detection

Existing object detection methods based on deep neural networks can be divided into anchor-based and anchor-free detectors. The anchor-based detectors use the preset boxes as anchors, which are trained to classify their categories and regress the offsets of coordinates. They can be further divided into multi-stage [9,36] and single-stage [35,29,8] detectors. As the representative multi-stage detector, Faster R-CNN [36] uses a region proposal network to generate proposals that probably contain objects and then predicts their categories and refines the proposals in the second stage. Considering the large computation cost of the multi-stage detectors, YOLO [35], as the representative single-stage detector, is proposed to use a fully convolutional network to predict both the bounding boxes and categories. It is further improved by applying feature pyramid [29], deconvolutional layers [8] and focal loss [27] to treat the various object scales, the semantic information of features, and the unbalance of positives and negatives, respectively. Many anchor-free detectors [40,49] are proposed to avoid empirically setting and tedious calculation of the anchors. Although applying deeper and wider networks can

often improve the performance of detectors, it is too computationally expensive in many resource-limited applications.

2.2 Knowledge Distillation

Knowledge distillation [16,47,14] is proposed by Hinton et al. [16] in the image classification task to transfer knowledge of a cumbersome teacher model into a compact student model. There are two main aspects of knowledge distillation: knowledge construction and knowledge transfer. For the first aspect, knowledge mainly consists of three types [10]: the feature-based knowledge, i.e., activations of intermediate feature [37,17,45,15], the relation-based knowledge, i.e., structures in the embedding space [34,31,39,42,25], and the response-based knowledge, i.e., the soft target of the output layers [16]. For the second aspect to effectively transfer the knowledge to the student. [16] applies a temperature factor to control the softness of the probability distribution over classes. [37] adds a regression layer as a bridge to match dimensions of the features. Such knowledge can be viewed as local knowledge since only the instance-level knowledge in the single feature space, e.g., the teacher’s is considered.

For distilling an object detector [3,11,48,46,23,22], more attention is paid on *constructing knowledge* due to the extreme imbalance over the foreground/background areas and the numbers of instances among different classes. [43] aims at keeping the balance between foreground and background features by distilling on the areas around ground-truth boxes, while [24] distills on high-level features within the equivalently sampled foreground and background proposals by referring to the ground-truth boxes. [4] is recently proposed to distill features in anchors where there are the most discrepancies of confidence between the student and the teacher model. [38] proposes to gradually reduce the distillation penalty to balance the two targets of detection and distillation. However, existing methods regard the activations or the relations between all instances as the local knowledge to distill object detectors, which suffers from the noises, e.g., the ambiguous instances or the detection errors from the teacher.

3 Method

In this section, we detail the proposed framework for distilling object detectors with global knowledge. As shown in Fig. 2, the overall framework consists of two modules: a prototype generation module (PGM) to find class-wise prototypes for bridging the two feature spaces, and a robust distilling module (RDM) to construct and distill the reliable global knowledge based on the prototypes.

3.1 Prototype Generation Module

The knowledge of a deep model can be viewed as the representation of its feature space, which can be approximated by a small set of basis vectors from the view of dictionary learning [41,20], as shown in Fig 1 (left). Concretely, let

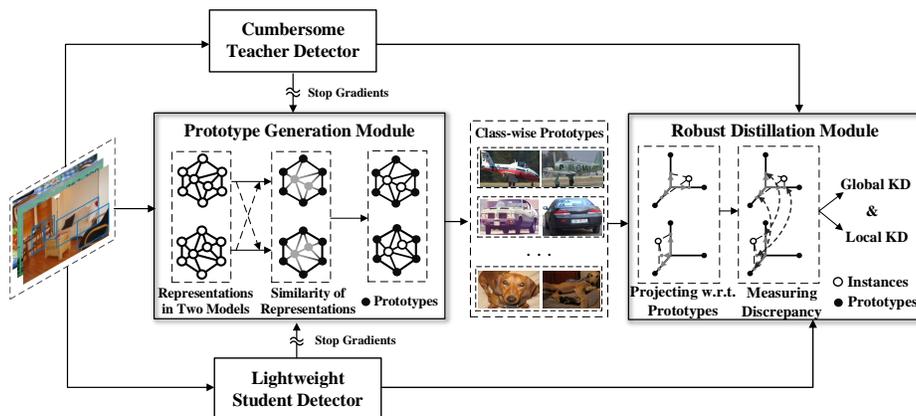


Fig. 2. The proposed framework for distilling object detectors with global knowledge. A prototype generation module (PGM) is first deployed to find the prototypes for each class based on the similarity of their representations in TS -space. A robust distillation module (RDM) is then designed to construct reliable global knowledge w.r.t. the prototypes and measure their discrepancy for robust knowledge distillation.

$\mathbf{F} = \{\mathbf{f}_i\}_{i=1}^N \in \mathbb{R}^{D \times N}$ be features of N instances in the feature space of D dimensions. The K ($K \ll N$) basis vectors $\mathbf{G} = \{\mathbf{g}_i\}_{i=1}^K \subset \mathbf{F}$ of a single feature space can be selected by minimizing the reconstruction errors of all instances:

$$\mathbf{G}, \mathbf{W} = \arg \min_{\mathbf{G}, \mathbf{W}} \|\mathbf{F} - \mathbf{G}\mathbf{W}\|_2^2 + \lambda \|\mathbf{W}\|_1^2, \quad (1)$$

where $\mathbf{W} \in \mathbb{R}^{K \times N}$ is the representation of all samples \mathbf{F} w.r.t. the basis vectors \mathbf{G} . The last regularized term weighted by λ helps to learn a sparse \mathbf{W} , which makes the basis vectors \mathbf{G} representative.

In the knowledge distillation task, there are two different feature spaces that are represented by the teacher and student detectors, namely TS -space. Thus, a more intrinsic approach is to find a group of common basis vectors in TS -space, which bridge the gap between the two spaces and reduce the difficulty of distillation. Following the above considerations, the prototype generation module (PGM) aims at finding K instances as the basis vectors, dubbed *prototypes*. In this way, other instances can be represented by prototypes with minimum reconstruction errors in each of the feature spaces. Meanwhile, the representing discrepancy based on the prototypes between two feature spaces should also be small such that it is easier to transform one feature space to another.

Let $\mathbf{F}_t = \{\mathbf{f}_i^t\}_{i=1}^N \in \mathbb{R}^{D_t \times N}$ and $\mathbf{F}_s = \{\mathbf{f}_i^s\}_{i=1}^N \in \mathbb{R}^{D_s \times N}$ be the N instances in the feature spaces of the teacher and the student detectors, respectively. D_t and D_s are the dimensions of the two feature spaces. The K prototypes can be grouped as $(\mathbf{G}_t, \mathbf{G}_s) = \{(\mathbf{g}_i^t, \mathbf{g}_i^s)\}_{i=1}^K$, where $\mathbf{G}_t = \mathbf{F}_{i(t)} \subset \mathbf{F}_t$ and $\mathbf{G}_s = \mathbf{F}_{i(s)} \subset \mathbf{F}_s$ are the subset of all instances in TS -space. $i(s)$, $i(t)$ are the indexing sets.

Notice that the prototypes are the common basis vectors of *TS-space*. Thus, they can be generated from all the instances by minimizing the reconstruction errors in both of the two feature spaces as well as a regularization of the representing consistency with a trade-off weight λ :

$$\begin{aligned} & \|\mathbf{F}_t - \mathbf{F}_{\mathbf{i}(t)} \mathbf{W}_t\|_2^2 + \|\mathbf{F}_s - \mathbf{F}_{\mathbf{i}(s)} \mathbf{W}_s\|_2^2 + \lambda \|\mathbf{W}_s - \mathbf{W}_t\|_2^2 \\ \text{s.t. } & \mathbf{i}(s) = \mathbf{i}(t) \quad \text{and} \quad |\mathbf{i}(s)| = |\mathbf{i}(t)| = K \end{aligned} \quad (2)$$

$\mathbf{W}_t = \{w_{j,i}^t\}_{K \times N}$ and $\mathbf{W}_s = \{w_{j,i}^s\}_{K \times N}$ are the representations, i.e., coordinates, of the N instances w.r.t. the K prototypes in the feature spaces of the teacher and the student. The last term in Eq. 2 requires the representations of instances in *TS-space* are similar such that the discrepancy of the relations is small, as illustrated in Fig. 1 (right). The constraint $\mathbf{i}(s) = \mathbf{i}(t)$ requires that a prototype is indeed one instance represented in two feature spaces, and total K prototypes are selected. In this way, representations of instances w.r.t. the prototypes can be regarded as approximations of the feature space, i.e., the global knowledge, as shown in Fig. 1 (left), which allows penalizing the difference of the relations between instances and prototypes in two feature spaces for knowledge transfer. Besides, the discrepancy of relations w.r.t. the prototypes can be used as the robustness cue for knowledge transfer, as illustrated in Fig. 1.

We show an approximate solution of the problem in Eq. 2 through a variant of matching pursuit [33], which is indeed a greedy algorithm yet very efficient. To select the $(n+1)^{\text{th}}$ prototype ($\mathbf{g}_{n+1}^t, \mathbf{g}_{n+1}^s$), we first define the residuals $\mathbf{r}_{n,i}^t$ and $\mathbf{r}_{n,i}^s$ w.r.t. the selected n prototypes as follows:

$$\mathbf{r}_{n,i}^t \triangleq \mathbf{f}_i^t - \sum_{k=1}^n \mathbf{g}_k^t w_{k,i}^t, \quad \mathbf{r}_{n,i}^s \triangleq \mathbf{f}_i^s - \sum_{k=1}^n \mathbf{g}_k^s w_{k,i}^s. \quad (3)$$

The objective in Eq. 2 w.r.t. the $(n+1)^{\text{th}}$ prototype can be written by

$$\mathcal{L}_{n+1} = \sum_{i=1}^N \|\mathbf{r}_{n+1,i}^t\|_2^2 + \sum_{i=1}^N \|\mathbf{r}_{n+1,i}^s\|_2^2 + \lambda \sum_{i=1}^N \sum_{k=1}^{n+1} (w_{k,i}^t - w_{k,i}^s)^2. \quad (4)$$

The optimal $w_{n+1,i}^t$ and $w_{n+1,i}^s$ can be obtained by making the derivative of the \mathcal{L}_{n+1} with respect of $w_{n+1,i}^t$ and $w_{n+1,i}^s$ to zero. Then, we have

$$w_{n+1,i}^t = \frac{\langle \mathbf{r}_{n,i}^t, \mathbf{g}_{n+1}^t \rangle + \lambda w_{n+1,i}^s}{\lambda + \|\mathbf{g}_{n+1}^t\|_2^2}, \quad w_{n+1,i}^s = \frac{\langle \mathbf{r}_{n,i}^s, \mathbf{g}_{n+1}^s \rangle + \lambda w_{n+1,i}^t}{\lambda + \|\mathbf{g}_{n+1}^s\|_2^2}. \quad (5)$$

We detail the derivation and show the closed-form solution of Eq. 5 in the supplemental materials, where we also show more analysis about the relationship between global knowledge and relation-based knowledge. The overall algorithm for generating prototypes is summarized in Alg. 1. Notice that we separately generate prototypes for each class.

Algorithm 1 Algorithm for selecting prototypes in PGM.

Input:
 $\{(\mathbf{f}_i^t, \mathbf{f}_i^s)\}_{i=1}^N$: features of N instances in TS -space;

Parameter:
 K : number of prototypes to be selected;

 λ : regularization weight;

Output:
 \mathcal{I} : index set of the prototypes

- 1: initialize $n = 0$, the residuals $\mathbf{r}_{0,i}^t = \mathbf{f}_i^t$, and $\mathbf{r}_{0,i}^s = \mathbf{f}_i^s \quad \forall i = 1, \dots, N$;
 - 2: **while** $n < K$ **do**
 - 3: compute the optimal $w_{n+1,i}^s$ and $w_{n+1,i}^t$ by Eq. 5;
 - 4: compute the \mathcal{L}_{n+1}^k with Eq. 4 for each instance by setting $\mathbf{g}_{n+1}^s = \mathbf{f}_k^s$ and $\mathbf{g}_{n+1}^t = \mathbf{f}_k^t, \forall k = 1, \dots, N$;
 - 5: append the index k^* into \mathcal{I} where $k^* = \arg \min_k \{\mathcal{L}_{n+1}^k\} \quad \forall (\mathbf{g}_k^s, \mathbf{g}_k^t) \in \{(\mathbf{f}_i^s, \mathbf{f}_i^t)\}_{i=1}^N$; set $\mathbf{g}_{n+1}^t = \mathbf{f}_{k^*}^t$ and $\mathbf{g}_{n+1}^s = \mathbf{f}_{k^*}^s$;
 - 6: update the residuals $\mathbf{r}_{n+1,i}^t$ and $\mathbf{r}_{n+1,i}^s$ by Eq. 3;
 - 7: set $n = n + 1$;
 - 8: **end while**
 - 9: **Return:** \mathcal{I}
-

3.2 Robust Distillation Module

In this section, we focus on *global knowledge construction* and *robust knowledge transferring* by a robust distillation module (RDM) based on the prototypes.

Identifying the knowledge. By referring to the prototypes, the global knowledge, i.e., the representations of instances on the common basis vectors in the two feature spaces, can be naturally constructed by measuring the representation between the instances and the prototypes.

Specifically, let the features of the j^{th} instance in the i^{th} image be $\mathbf{f}_{i,j}^t$ and $\mathbf{f}_{i,j}^s$ in the feature spaces of the teacher and the student detectors, respectively. For an instance with a pair of features $(\mathbf{f}_{i,j}^t, \mathbf{f}_{i,j}^s)$ in TS -space, they can be separately projected onto the common basis vectors $(\mathbf{G}_t, \mathbf{G}_s)$ in each space as $\mathbf{A}_{i,j}^t = \mathcal{P}_{\mathbf{G}_t}(\mathbf{f}_{i,j}^t)$ and $\mathbf{A}_{i,j}^s = \mathcal{P}_{\mathbf{G}_s}(\mathbf{f}_{i,j}^s)$. \mathcal{P} is the projection function. The project coefficients $\mathbf{A}_{i,j}^t$ and $\mathbf{A}_{i,j}^s$ can be calculated exactly the same as in Eq. 5.

Thus, the global knowledge can be transferred by minimizing:

$$\mathcal{L}_{\text{global}} = \frac{1}{2NK} \sum_{i=1}^n \sum_{j=1}^{n_i} \sigma_{i,j} \|\mathbf{A}_{i,j}^s - \mathbf{A}_{i,j}^t\|_2^2, \quad (6)$$

where $N = \sum_{i=1}^n n_i$ are the total number of instances. n is the number of images and n_i is the number of instances in the i -th image. $\sigma_{i,j}$ is the weight that reveals how reliable the knowledge is and will be discussed later.

For the local feature-based knowledge, we follow [43] identifying the knowledge as the features of the regions that overlap with any ground-truth boxes larger

Algorithm 2 The proposed knowledge distilling process.

Input: teacher detector \mathcal{T} , student detector \mathcal{S} , prototype updating period T and maximum training epochs T_m .

- 1: let e be the current training epoch and set $e = 0$;
 - 2: **while** $e < T_m$ **do**
 - 3: **if** $\text{mod}(e, T) == 0$ **then**
 - 4: extract features of instances \mathbf{F}_t and \mathbf{F}_s from the teacher \mathcal{T} and current student (at e -th epoch) \mathcal{S}^e , respectively;
 - 5: updating and bootstrapping prototypes $(\mathbf{G}_t, \mathbf{G}_s)$ for each class by minimizing Eq. 2 based on \mathcal{T} and \mathcal{S}^e (see Alg. 1);
 - 6: **end if**
 - 7: training the student detector for one epoch by minimizing Eq. 9;
 - 8: set $e = e + 1$;
 - 9: **end while**
-

than a threshold. Thus, the local feature-based knowledge can be defined as:

$$\mathcal{L}_{\text{local}}^{\text{feat}} = \frac{1}{2N} \sum_{i=1}^n \sum_{j=1}^{n_i} \sigma_{i,j} \|\mathcal{H}(\mathbf{f}_{i,j}^s) - \mathbf{f}_{i,j}^t\|_2^2, \quad (7)$$

where \mathcal{H} is an adaptation function, e.g., a 1×1 convolutional layer with ReLU activation in our paper, that transforms the features of the student into the feature space of the same dimensions as the teacher’s.

For the local response-based knowledge, we use the proposals and apply the RoI-align [36] to get the prediction inside the regions. The KL-divergence weighted by $\sigma_{i,j}$ is used on the predicting logits between the teacher and the student, and denoted as $\mathcal{L}_{\text{local}}^{\text{resp}}$.

Robustly distilling the knowledge. Since the knowledge from the teacher might be noisy, especially on ambiguous instances, a robust knowledge transferring approach is required to distinguish noisy knowledge and mitigate transferring them to the student. Inspired from co-teaching [18,32,12] to alleviate the noise from multiple views, the student might also have a voice in discriminating the noisy knowledge. Based on the observations that reliable knowledge should have similar representations under the measurement from two models, shown in Fig. 1, the robustness of knowledge can be estimated by the discrepancy of representations in *TS-space*. Thus, the weight $\sigma_{i,j}$ for fine-grained knowledge distillation can be approximated as

$$\sigma_{i,j} = 1 - \|\mathbf{A}_{i,j}^s - \mathbf{A}_{i,j}^t\|_2. \quad (8)$$

$\sigma_{i,j}$ describes the similarity of the representations between the instance and the prototypes in the *TS-space*. It is indeed heavily related to the last term in Eq. 2, where we concentrate more on the instances with small discrepancy w.r.t. the prototypes for both global and local knowledge transfer.

Table 1. Knowledge distillation results on COCO dataset with different detectors. Some results are missing since we cannot find the performance report in their papers.

Method	mAP	AP ₅₀	AP ₇₅	AP _s	AP _m	AP _l	mAR	AR _s	AR _m	AR _l
Faster-Res101 (teacher)	39.8	60.1	43.3	22.5	43.6	52.8	53.0	32.8	56.9	68.6
Faster-Res50 (student)	38.4	59.0	42.0	21.5	42.1	50.3	52.0	32.6	55.8	66.1
FGFI [43]	39.3	59.8	42.9	22.5	42.3	52.2	52.4	32.2	55.7	67.9
DeFeat [11]	40.3	60.9	44.0	23.1	44.1	53.4	53.7	33.3	57.7	69.1
FBKD [46]	40.2	60.4	43.6	22.8	43.8	53.2	53.4	32.7	57.1	68.8
GID [4]	40.2	60.8	43.6	23.6	43.9	53.0	53.7	33.6	57.7	68.6
Ours	40.6	61.0	44.0	23.4	44.4	53.3	53.8	33.9	57.9	69.2
Retina-Res101 (teacher)	38.9	58.0	41.5	21.0	42.8	52.4	54.8	33.4	59.3	71.2
Retina-Res50 (student)	37.4	56.7	39.6	20.0	40.7	49.7	53.9	33.1	57.7	70.2
FGFI [43]	38.6	58.7	41.3	21.4	42.5	51.5	54.6	34.7	58.2	70.4
GID [4]	39.1	59.0	42.3	22.8	43.1	52.3	55.3	36.7	59.1	71.1
DeFeat [11]	39.3	58.2	42.1	21.7	42.9	52.9	55.1	33.9	59.6	71.5
FBKD [46]	39.3	58.8	42.0	21.2	43.2	53.0	55.4	34.6	59.7	72.2
FR [5]	39.3	58.8	42.0	21.5	43.3	52.6	-	-	-	-
PFI [23]	39.6	-	-	21.4	44.0	52.5	-	-	-	-
Ours	39.8	58.6	42.6	21.8	43.5	53.5	55.8	34.1	60.0	72.2

3.3 Optimization

The overall objective for distilling object detectors can be summarized as:

$$\mathcal{L}_{\text{kd}} = \mathcal{L}_{\text{det}} + \alpha_1 \mathcal{L}_{\text{global}} + \alpha_2 \mathcal{L}_{\text{local}}^{\text{feat}} + \alpha_3 \mathcal{L}_{\text{local}}^{\text{resp}}, \quad (9)$$

where \mathcal{L}_{det} is the original detection objective defined by the student detector. α_1 , α_2 , and α_3 weigh the global and local knowledge transfer. For detectors with FPN [26] using multiple feature maps for prediction, we independently apply the PGM and the RDM on each of the feature maps. Since the student is gradually optimized and the relations are changed, the prototypes should be updated when training the student. For efficiency, the prototypes are bootstrapped and updated every T epochs. Both the student and the teacher detectors are pre-trained on the task-relevant dataset to extract features of instances and generate the prototypes. The overall proposed distilling algorithm is summarized in Alg. 2

4 Experiments

We perform experiments with the representative single-stage and two-stage detectors, namely, RetinaNet [27] and Faster R-CNN [36] on the PASCAL VOC [7] and COCO [28] detection benchmarks. We follow the common settings that use both VOC 07 and 12 *trainval* split for training and VOC 07 *test* split for test. For the COCO dataset, the *train* split are used for training while the *val* split are used for test. Unless otherwise specified, the hyper-parameters are set as $K = 10$, $\lambda = 10$ and $T = 1$. The distilling weights α_1 , α_2 , and α_3 are set to

1.0, 1.0, 5.0, respectively. The student detector is trained through $2\times$ learning schedule on 8 Tesla V100 32G GPUs. The input images are resized as large as 1333×800 while keeping the aspect ratio. Other standard augmentations, e.g., the photometric distortion, are applied as the settings in MMDetection [19]. The ResNet101 and ResNet50 [13] backbones are used for the teacher and the student detectors, respectively. We also validate our methods with larger teachers, e.g., Cascade Mask R-CNN [2] with ResNext-101 [44]. More implementation details are included in the supplemental material.

4.1 Comparison with existing methods on VOC and COCO datasets

We first evaluate our method on VOC and COCO datasets with the representative two-stage detector (Faster R-CNN) and single-stage detector (RetinaNet). As shown in Table 1 and Table 2, all student models are significantly improved by our knowledge distillation algorithm, e.g., 2.2% and 2.4% mAP on the COCO dataset and 2.5% and 2.5% mAP on the VOC dataset for both detectors. Moreover, they even surpass the teacher detector within a large margin, e.g., 0.8%, 0.9% on COCO dataset for both detectors. As we form the global knowledge as the ensemble of both the student and the teacher detectors and use common basis vectors to bridge the two feature spaces for distilling, the proposed method shows more potential to achieve a further gain compared to the teacher detectors.

We also compare our method with the SOTA detection distillation methods with the same teacher and student detectors. Table 1 and Table 2 show that the proposed method achieves best mAP on COCO and VOC datasets. Notice that GID [4] applies all the three types of local knowledge, i.e., feature-based, relation-based, and response-based knowledge for distilling. However, the proposed method shows further improvement on both COCO and VOC datasets, e.g, 0.7% mAP gain for distilling the RetinaNet. It reveals that distilling the knowledge by forcing the student to absolutely behave the same as the teacher still leads to sub-optimal since the local knowledge represented by the ambiguous instances is hard to transfer and will hurt the distillation. The proposed method shows a more promising way by looking for a group of common basis vectors, i.e., the prototypes, for bridging the gap of the two feature spaces and forming as well as distilling the global knowledge based on the prototypes in a more robust way. Moreover, the results in Table 1 and Table 2 demonstrate that our method is capable to be applied to various detection frameworks.

4.2 Effects of the prototypes in robust knowledge distillation

To verify the advantages of the prototypes bridging the two feature spaces for global and local knowledge distillation, we conduct ablation experiments on the VOC dataset with $1\times$ learning schedule. ResNet101-based and ResNet50-based Faster R-CNN are used as the teacher and the student detectors, respectively.

We first separately apply $\mathcal{L}_{\text{global}}$, $\mathcal{L}_{\text{local}}^{\text{feat}}$ and $\mathcal{L}_{\text{local}}^{\text{resp}}$ in Eq. 9 for knowledge distillation. In our framework, the global knowledge is formed as the projections of instances w.r.t. the prototypes, while the feature-based and response-based

Table 2. Knowledge distillation results on Pascal VOC dataset with different detectors.

Method	Faster R-CNN Res101-50			RetinaNet Res101-50		
	mAP	AP ₅₀	AP ₇₅	mAP	AP ₅₀	AP ₇₅
teacher	56.3	82.7	62.6	58.2	82.0	63.0
student	54.2	82.1	59.9	56.1	80.9	60.7
FitNet [37]	55.0	82.2	61.2	56.4	81.7	61.7
FGFI [43]	55.3	82.1	61.1	55.6	81.4	60.5
FBKD [46]	55.4	82.0	61.3	56.7	81.9	61.9
ICD [22]	56.4	82.4	63.4	57.7	82.4	63.5
GID [4]	56.5	82.6	61.6	57.9	82.0	63.2
Ours	56.7	82.9	61.9	58.6	82.4	64.2

Table 3. Ablation experiment on separately applying the global knowledge $\mathcal{L}_{\text{global}}$, feature-based local knowledge $\mathcal{L}_{\text{local}}^{\text{feat}}$ and response-based local knowledge $\mathcal{L}_{\text{local}}^{\text{resp}}$ in Eq. 9 on VOC dataset with $1\times$ learning schedule.

Module	Student	Faster R-CNN Res101-50					
$\mathcal{L}_{\text{local}}^{\text{feat}}$		✓			✓		✓
$\mathcal{L}_{\text{global}}$			✓			✓	✓
$\mathcal{L}_{\text{local}}^{\text{resp}}$				✓	✓	✓	✓
AP ₅₀	81.3	82.0	82.4	82.2	82.6	82.4	82.9

local knowledge is weighted through the discrepancy of the projections. Table 3 shows that the prototypes can boost the global and local knowledge distillation by a large margin. By applying both the global and local knowledge, we show 1.6% performance gain compared to the student detector with the $1\times$ learning schedule, which also surpasses the teacher detector with the mAP 82.4%.

Furthermore, we also extend some existing methods based on the prototypes. DeFeat [11], RKD [34] and Vanilla-KD [16] are the representative feature-based, relation-based and response-based local knowledge distillation methods. We directly use the released source codes of FBKD and carefully re-implement the RKD (as RKD[†]) and Vanilla-KD as (Vanilla-KD[†]) for distilling the object detectors. Then, we apply the prototypes separately: as for RKD [34], we form the global knowledge by projecting the instances w.r.t. the prototypes; as for DeFeat [11] and Vanilla-KD [16], we apply the distilling weight defined in Eq. 8, which is measured by the discrepancy w.r.t. the prototypes. Table 4 shows the consistent performance gain among those three knowledge distillation methods with the prototypes, which shows the effectiveness of prototypes for both constructing more reliable global knowledge and more robust knowledge transfer. The implementation details by combining prototypes with those distilling methods are included in the supplemental materials.

Table 4. Ablation experiment by combining the prototypes with the existing representative methods for feature-based [11], relation-based [34], and response-based [16] local knowledge distillation, respectively.

Method	Student	DeFeat [11]	RKD [†] [34]	Vanilla-KD [†] [16]
+prototypes		✓	✓	✓
AP₅₀	81.3	82.0 82.4	81.6 82.0	81.8 82.2

Table 5. Ablation experiments on the hyperparameters α_1 , α_2 , α_3 , λ , K , and T .

α_1	0.1	0.5	1.0	1.2	α_2	0.5	0.8	1.0	1.5
AP ₅₀	82.1	82.3	82.9	82.6	AP ₅₀	82.3	82.4	82.9	82.6
α_3	1	5	10	20	λ	1	10	50	100
AP ₅₀	82.3	82.9	82.3	82.2	AP ₅₀	82.3	82.9	82.6	82.0
K	1	5	10	20	T	0.5	1	2	3
AP ₅₀	82.0	82.3	82.9	82.5	AP ₅₀	82.8	82.9	82.3	82.2

4.3 Analysis of the hyperparameters

We investigate the updating periods T in Alg. 2 of the prototypes for knowledge distillation on the VOC dataset. Since the student detector is updated during training, the prototypes and their features \mathbf{G}_s should be updated. Table 5 shows that as the period T increasing, the performance slightly decreases. It is because the bootstrapped prototypes are approximations of the basis vectors of updated student detector, which results in some bias when forming the global knowledge as well as computing the discrepancy in Eq. 8. Table 5 also shows ablation experiments on the three weights α_1 , α_2 , and α_3 of the three terms in Eq. 9, the number of selected prototypes K and the similarity regularization weight λ in Eq. 2. The results in Table 5 show that the proposed method is relatively robust to the hyperparameters, which achieves better performance than the student in a wide range of hyperparameters.

4.4 Analysis on the prototype generation methods

In our framework, the prototypes play roles as the common basis vectors in both the feature spaces of the teacher and the student. They are selected by minimizing the reconstruction errors among instances in TS -space as defined in Eq. 2. We also compare the proposed prototype generation algorithm in Alg. 1 with some other similar methods, e.g., the K-means and the DBSCAN [6]. Besides, we also deliberately select the same number of ambiguous instances, e.g., small or truncated instances, as the prototypes for comparison. In Table 6, we show the performance of knowledge distillation based on those prototype generation methods. We find that the cluster-like algorithms, e.g., the K-means or the DBSCAN [6], fail to improve the distillation by comparing the results in Table 3, because those algorithms are applied only in the single feature space

Table 6. Ablation experiments on the effect of different prototype generation methods for knowledge distillation. For the cluster-like algorithms, e.g., K-Means and DBSCAN [6], we apply them separately on the feature space of either the teacher or the student.

Method	K-Means		DBSCAN [6]		Ambiguous	Ours
Features	Student	Teacher	Student	Teacher	-	Both
AP₅₀	82.3	82.1	82.4	82.2	81.8	82.9

and can hardly bridge the two feature spaces of the teacher and the student. The poor performance by selecting the ambiguous instances as the prototypes further verify the importance of selecting the representative instances as the prototypes. Otherwise, it will bring large discrepancy as shown in Fig. 1 (right), and increase the difficulty of knowledge distillation.

4.5 Distilling with larger teacher

The larger teacher will achieve better performance, which might also bring an extra bonus for knowledge distillation. Following the common settings with the existing methods [11,43,38], on the VOC dataset, we use the Faster R-CNN with the backbones ResNet152 and ResNet50 as the teacher and the student. For a fair comparison, we follow DeFeat [11] by using $1\times$ learning schedule. On the COCO dataset, we follow FBKD [46] by applying ResNeXt101-based [44] Cascade Mask R-CNN [2] as the teacher detector and the ResNet50-based Faster R-CNN as the student. The $2\times$ learning schedule is used as in FBKD [46]. In Table 8 and Table 7, we show the performance of knowledge distillation with larger teachers on VOC and COCO datasets, respectively. The proposed method can still achieve the best performance on the VOC dataset, with the 0.6% mAP advantage w.r.t. DeFeat [11]. On the COCO dataset, we achieve comparable performance as the FBKD [46] with the much larger teacher and heterogeneous backbone. The performance with larger teachers further shows the proposed method can be applied in various detection frameworks with the same hyperparameters.

4.6 Analysis of noisy knowledge transferring

In Figure 3, we also illustrate some wrong detection in red boxes, *e.g.*, false positives and inaccurately located instances, from the teacher detector that are

Table 7. Knowledge distillation results with larger teacher on COCO dataset.

Method	Backbone	mAP	AP ₅₀	AP ₇₅	AP _s	AP _m	AP _l
Cascade R-CNN (teacher)	ResNext101	47.3	66.3	51.7	28.2	51.7	62.7
Faster R-CNN (student)	ResNet50	38.4	59.0	42.0	21.5	42.1	50.3
FBKD [46]	ResNet50	41.5	62.2	45.1	23.5	45.0	55.3
Ours	ResNet50	41.5	61.9	45.1	23.5	45.1	55.4



Fig. 3. Illustration of detection results via different knowledge distillation methods, e.g., FGFI [43], our re-implemented RKD [34], DeFeat [11], and ours. Some noisy knowledge of the teacher are transferred to the student (marked in red boxes). Best view in color.

Table 8. Knowledge distillation results with larger teacher on VOC dataset.

Method	Teacher	Student	Faster R-CNN ResNet152-50			
			FGFI[43]	TADF[38]	DeFeat[11]	Ours
AP_{50}	83.1	81.3	81.6	81.7	82.3	82.9

transferred to the student. Our method shows more promising results against noisy knowledge transferring and is capable to surpass the performance of the teacher detector. More quantitative analysis is discussed in the supplementary.

5 Conclusion

In this paper, we propose a novel knowledge distillation framework with global knowledge. The prototype generation module is first designed to find a group of common basis vectors, i.e., the *prototypes*, by minimizing the reconstruction errors in both the feature spaces of the teacher and the student. The robust distillation module is then applied to (1) construct the global knowledge by projecting the instances w.r.t. the prototypes, and (2) robustly distill the global and local knowledge by measuring their discrepancy in the two spaces. Experiments show that the proposed method achieves state-of-the-art performance on two popular detection frameworks and benchmarks. The extensive experimental results show that the proposed method can be easily stretched with larger teachers and the existing knowledge distillation methods to obtain further improvement.

References

1. Bucila, C., Caruana, R., Niculescu-Mizil, A.: Model compression. In: SIGKDD. pp. 535–541 (2006)
2. Cai, Z., Vasconcelos, N.: Cascade R-CNN: high quality object detection and instance segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.* **43**(5), 1483–1498 (2021)
3. Chen, G., Choi, W., Yu, X., Han, T.X., Chandraker, M.: Learning efficient object detection models with knowledge distillation. In: NeurIPS. pp. 742–751 (2017)
4. Dai, X., Jiang, Z., Wu, Z., Bao, Y., Wang, Z., Liu, S., Zhou, E.: General instance distillation for object detection. *CoRR* [abs/2103.02340](#) (2021)
5. Du, Z., Zhang, R., Chang, M., Zhang, X., Liu, S., Chen, T., Chen, Y.: Distilling object detectors with feature richness. *CoRR* [abs/2111.00674](#) (2021)
6. Ester, M., Kriegel, H., Sander, J., Xu, X.: A density-based algorithm for discovering clusters in large spatial databases with noise. In: Simoudis, E., Han, J., Fayyad, U.M. (eds.) *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96)*, Portland, Oregon, USA. pp. 226–231. AAAI Press (1996)
7. Everingham, M., Gool, L.V., Williams, C.K.I., Winn, J.M., Zisserman, A.: The pascal visual object classes (VOC) challenge. *IJCV* **88**(2), 303–338 (2010)
8. Fu, C., Liu, W., Ranga, A., Tyagi, A., Berg, A.C.: DSSD : Deconvolutional single shot detector. *CoRR* [abs/1701.06659](#) (2017)
9. Girshick, R.: Fast R-CNN. In: ICCV. pp. 1440–1448 (2015)
10. Gou, J., Yu, B., Maybank, S.J., Tao, D.: Knowledge distillation: A survey. *CoRR* [abs/2006.05525](#) (2020)
11. Guo, J., Han, K., Wang, Y., Wu, H., Chen, X., Xu, C., Xu, C.: Distilling object detectors via decoupled features. *CoRR* [abs/2103.14475](#) (2021)
12. Han, B., Yao, Q., Yu, X., Niu, G., Xu, M., Hu, W., Tsang, I.W., Sugiyama, M.: Co-teaching: Robust training of deep neural networks with extremely noisy labels. In: NeurIPS. pp. 8536–8546 (2018)
13. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR. pp. 770–778. CS (2016)
14. Heo, B., Kim, J., Yun, S., Park, H., Kwak, N., Choi, J.Y.: A comprehensive overhaul of feature distillation. In: ICCV. pp. 1921–1930 (2019)
15. Heo, B., Lee, M., Yun, S., Choi, J.Y.: Knowledge transfer via distillation of activation boundaries formed by hidden neurons. In: AAA. pp. 3779–3787 (2019)
16. Hinton, G.E., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. *CoRR* [abs/1503.02531](#) (2015)
17. Huang, Z., Wang, N.: Like what you like: Knowledge distill via neuron selectivity transfer. *CoRR* [abs/1707.01219](#) (2017)
18. Jiang, L., Zhou, Z., Leung, T., Li, L., Fei-Fei, L.: Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels. In: ICML. vol. 80, pp. 2309–2318 (2018)
19. Kai Chen, e.a.: Mmdetection: Open mmlab detection toolbox and benchmark. *CoRR* [abs/1906.07155](#) (2019)
20. Kreutz-Delgado, K., Murray, J.F., Rao, B.D., Engan, K., Lee, T., Sejnowski, T.J.: Dictionary learning algorithms for sparse representation. *Neural Comput.* **15**(2), 349–396 (2003)
21. Le, E., Kokkinos, I., Mitra, N.J.: Going deeper with lean point networks. In: CVPR. pp. 9500–9509 (2020)

22. Li, G., Li, X., Wang, Y., Zhang, S., Wu, Y., Liang, D.: Knowledge distillation for object detection via rank mimicking and prediction-guided feature imitation. *CoRR abs/2112.04840* (2021)
23. Li, G., Li, X., Wang, Y., Zhang, S., Wu, Y., Liang, D.: Knowledge distillation for object detection via rank mimicking and prediction-guided feature imitation. *CoRR abs/2112.04840* (2021)
24. Li, Q., Jin, S., Yan, J.: Mimicking very efficient network for object detection. In: *CVPR*. pp. 7341–7349 (2017)
25. Li, X., Wu, J., Fang, H., Liao, Y., Wang, F., Qian, C.: Local correlation consistency for knowledge distillation. In: *ECCV*. vol. 12357, pp. 18–33 (2020)
26. Lin, T., Dollár, P., Girshick, R.B., He, K., Hariharan, B., Belongie, S.J.: Feature pyramid networks for object detection. In: *CVPR*. pp. 936–944 (2017)
27. Lin, T., Goyal, P., Girshick, R.B., He, K., Dollár, P.: Focal loss for dense object detection. In: *ICCV*. pp. 2999–3007 (2017)
28. Lin, T., Maire, M., Belongie, S.J., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft COCO: common objects in context. In: *ECCV*. vol. 8693, pp. 740–755 (2014)
29. Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S.E., Fu, C., Berg, A.C.: SSD: single shot multibox detector. In: *ECCV*. vol. 9905, pp. 21–37 (2016)
30. Liu, Y., Cao, J., Li, B., Yuan, C., Hu, W., Li, Y., Duan, Y.: Knowledge distillation via instance relationship graph. In: *CVPR*. pp. 7096–7104 (2019)
31. Liu, Y., Cao, J., Li, B., Yuan, C., Hu, W., Li, Y., Duan, Y.: Knowledge distillation via instance relationship graph. In: *CVPR*. pp. 7096–7104 (2019)
32. Malach, E., Shalev-Shwartz, S.: Decoupling ”when to update” from ”how to update”. In: *NeurIPS*. pp. 960–970 (2017)
33. Mallat, S., Zhang, Z.: Matching pursuits with time-frequency dictionaries. *TIP* **41**(12), 3397–3415 (1993)
34. Park, W., Kim, D., Lu, Y., Cho, M.: Relational knowledge distillation. In: *CVPR*. pp. 3967–3976 (2019)
35. Redmon, J., Divvala, S.K., Girshick, R.B., Farhadi, A.: You only look once: Unified, real-time object detection. In: *CVPR*. pp. 779–788 (2016)
36. Ren, S., He, K., Girshick, R.B., Sun, J.: Faster R-CNN: towards real-time object detection with region proposal networks. In: *NeurIPS*. pp. 91–99 (2015)
37. Romero, A., Ballas, N., Kahou, S.E., Chassang, A., Gatta, C., Bengio, Y.: Fitnets: Hints for thin deep nets. In: *ICLR* (2015)
38. Sun, R., Tang, F., Zhang, X., Xiong, H., Tian, Q.: Distilling object detectors with task adaptive regularization. *CoRR abs/2006.13108* (2020)
39. Tian, Y., Krishnan, D., Isola, P.: Contrastive representation distillation. In: *ICLR* (2020)
40. Tian, Z., Shen, C., Chen, H., He, T.: FCOS: fully convolutional one-stage object detection. In: *ICCV*. pp. 9626–9635 (2019)
41. Tosic, I., Frossard, P.: Dictionary learning. *SPM* (2011)
42. Tung, F., Mori, G.: Similarity-preserving knowledge distillation. In: *ICCV*. pp. 1365–1374 (2019)
43. Wang, T., Yuan, L., Zhang, X., Feng, J.: Distilling object detectors with fine-grained feature imitation. In: *CVPR*. pp. 4933–4942 (2019)
44. Xie, S., Girshick, R.B., Dollár, P., Tu, Z., He, K.: Aggregated residual transformations for deep neural networks. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21–26, 2017*. pp. 5987–5995. IEEE Computer Society (2017)

45. Yim, J., Joo, D., Bae, J., Kim, J.: A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In: CVPR. pp. 7130–7138 (2017)
46. Zhang, L., Ma, K.: Improve object detection with feature-based knowledge distillation: Towards accurate and efficient detectors. In: ICLR (2021)
47. Zhang, Y., Lan, Z., Dai, Y., Zeng, F., Bai, Y., Chang, J., Wei, Y.: Prime-aware adaptive distillation. In: ECCV. vol. 12364, pp. 658–674 (2020)
48. Zheng, Z., Ye, R., Wang, P., Wang, J., Ren, D., Zuo, W.: Localization distillation for object detection. CoRR **abs/2102.12252** (2021)
49. Zhu, C., He, Y., Savvides, M.: Feature selective anchor-free module for single-shot object detection. In: CVPR. pp. 840–849 (2019)