# Dynamic Supervisor for Cross-dataset Object Detection

Ze Chen[a], Zhihang Fu[b], Jianqiang Huang[b], Mingyuan Tao[b], Shengyu Li[b], Rongxin Jiang[a], Xiang Tian[a], Yaowu Chen[a,*], Xian-Sheng Hua[b]

[a]*Zhejiang University, No.38 Zheda Road, Hangzhou, 310027, Zhejiang, China*
[b]*Alibaba Group, Hangzhou, 311121, Zhejiang, China*

## Abstract

The application of cross-dataset training in object detection tasks is complicated because the inconsistency in the category range across datasets transforms fully supervised learning into semi-supervised learning. To address this problem, recent studies focus on the generation of high-quality missing annotations. In this study, we first point out that it is not enough to generate high-quality annotations using a single model, which only looks once for annotations. Through detailed experimental analyses, we further conclude that hard-label training is conducive to generating high-recall annotations, while soft-label training tends to obtain high-precision annotations. Inspired by the aspects mentioned above, we propose a dynamic supervisor framework that updates the annotations multiple times through multiple-updated submodels trained using hard and soft labels. In the final generated annotations, both recall and precision improve significantly through the integration of hard-label training with soft-label training. Extensive experiments conducted on various dataset combination settings support our analyses and demonstrate the superior performance of the proposed dynamic supervisor.

*Keywords:* Cross-dataset object detection, Hard-label training, Soft-label training, Dynamic ensembling

## 1. Introduction

Fully supervised learning [1] has dominated the field of computer vision for decades. One of the most distinctive features of fully supervised learning is "data-driven learning" in which a certain annotated dataset [2, 3, 4] defines the capability

---

*Corresponding author
  *Email address:* cyw@mail.bme.zju.edu.cn (Yaowu Chen)

boundary of the trained model (the categories that are identifiable and those that are not). As a result, if one intends to expand the capability boundary, the training data must be significantly augmented with a larger category set, and this process is time-consuming and labor-intensive. Therefore, cross-dataset training [5, 6, 7, 8] has attracted the attention of academics seeking to avoid the unignorable costs of establishing new datasets. In cross-dataset training, several existing datasets only need to be merged and trained to expand the capability boundary to the union of their category sets without any extra image labeling costs.

Compared to classification [5] and image retrieval [6] tasks, cross-dataset training on object detection is much more complicated. This is because of the multi-label and multi-instance attributes of object detection. In an image, detection annotation always contains multiple categories and instances. A naive combination across detection datasets results in the incompleteness of the entire set. For example, the category *car*, which is annotated on dataset Alpha, is essentially not annotated at all on datasets Beta and Gamma. These incomplete annotation sets further transform fully supervised learning into weakly [9, 10, 11, 12] or semi-supervised learning [13, 14, 15]. As a result of the annotation insufficiency in each image, cross-dataset training encounters the problem of confusing positive and negative samples during sample selection.

Recent studies focus on generating missing annotations through detection submodels trained using each dataset individually to circumvent this difficulty. Continuing with the previous example, in the studies conducted by [16, 8], the missing "*car*" annotation is predicted on datasets Beta and Gamma through the submodel trained using dataset Alpha. We regard the approach employed in these studies as a static supervisor framework because *the submodel only looks once.* We can only hope that the submodel is strong enough to generate optimal annotations, i.e., both high recall and high precision. However, in practice, owing to the capability discrepancy across different models and categories, the quality of annotation generation cannot be guaranteed. There is no doubt that cross-dataset training for object detection still requires further exploration.

In this study, we first discuss the difference in the quality of annotation generation between hard-label training and soft-label training. Although the detection performance of both approaches is similar, the annotations generated through hard-label training tend to have higher confidence scores and more false positives than those generated through soft-label training. Therefore, we conclude that hard-label training is conducive to generating high-recall annotations, whereas soft-label training obtains high-precision annotations. Inspired by the above, we propose a dynamic supervisor framework for cross-dataset training. It leverages both the submodel from hard-label

training and that from soft-label training to expand and shrink the generated annotations dynamically and obtain a final annotation set with high recall and high precision. In other words, the proposed dynamic supervisor updates the annotations multiple times according to the multiple-updated submodels, avoiding the dilemma of *only looking once*. As shown in Figure 1, the annotations (highlighted in blue) generated using the initial submodel only cover a portion of the *car* instances. The hard-label-training submodel then adds additional predictions (highlighted in yellow) to increase the recall rate. Afterward, the soft-label-training submodel re-checks the annotations and reduces the unreliable predictions (the blue dotted rectangle) to increase the precision rate. In this study, we reveal the implicit connections between hard- and soft-label training and the methods for pseudo-annotation ensembling. The dynamic supervisor framework is designed to utilize these connections and achieve superior performance. Extensive experiments conducted on three combinations of several existing datasets [17, 18, 19, 20] support our analyses and proposed method. The main contributions of this study are as follows:

1. We show that hard-label training and soft-label training are conducive to improving the recall and precision of predictions, respectively, and that their integration brings improvements to both.
2. We propose a dynamic supervisor framework, which dynamically polishes the annotations and adaptively selects predictions based on category.
3. We conduct detailed experiments to demonstrate the effectiveness of the proposed framework, and we achieve state-of-the-art performance on all three cross-dataset settings.

## 2. Related Work

Over the past decade, cross-dataset learning has become increasingly popular. In object detection, cross-dataset learning is significantly different from image classification or image retrieval tasks. When multiple datasets annotated with different category sets are merged, the incomplete annotation transforms fully supervised learning into semi-supervised learning [21, 22, 23]. Radosavovic et al. [21] proposed a method called data distillation, which ensembles predictions from multiple transformations of unlabeled data, using a single model to automatically generate new training annotations. Jeong et al. [22] proposed consistency constraints to enhance detection performance. Wu et al. [24] proposed a soft sampling method that re-weighs the gradients of RoIs as a function of overlap with positive instances. This method ensures that uncertain background regions are given a smaller weight compared to that of the hard-negatives when there are some unlabeled object instances in training images.

3

Yang et al. [25] proposed treating object detection as a positive-unlabeled problem, which removes the assumption that unlabeled regions must be negative. Chadwick and Newman [26] examined the effect of different types of label noise on the performance of an object detector and applied the co-teaching framework to improve the performance of the detector trained on a noisy dataset. These works exploited their methods to make full use of the unlabeled or noisy-labeled data, but could not expand the capability boundary of detection models.

Many existing works [8, 16, 27] attempt to tackle the cross-dataset object detection with pseudo-labeling methods. Abbasi et al. [27] proposed a computationally efficient self-supervised framework to create pseudo-labels for the unlabeled positive instances in the merged dataset in order to train the object detector jointly on both ground truths and pseudo labels. Rame et al. [16] proposed the use of models trained using several datasets individually to generate pseudo-annotations on other datasets, and they proposed a new classification loss called SoftSig to handle unreliable pseudo-annotations. Zhao et al. [8] exploited a pseudo-labeling approach and proposed loss functions that carefully integrated partial but correct annotations with complementary but noisy pseudo-labels.

In the training of supervised neural networks, the different usage of labels results in different properties of neural networks. Soft-label training can enhance the smoothness of output probabilities and prevent overconfident predictions [28]. Müller et al. [29] found that training a network using hard labels typically results in the correct logit being much larger than any of the incorrect logits, and it also allows the incorrect logits to be significantly different from one another. Lukasik et al. [30] reported that label smoothing is effective as a technique for coping with label noise, and it improves the accuracy of both the clean and noisy parts of the data.

## 3. Analysis of Hard- and Soft-Label Training

We dedicate a separate section to dissect the differences in the quality of annotation generation between hard-label training and soft-label training. Without loss of generality, we adopt two datasets with different category sets for cross-dataset training. For the experiments in this section, we choose the single-shot detector (SSD) [31] to avoid disturbing factors introduced from complex detection frameworks.

### 3.1. Experiment Setting

**Dataset.** We sample images from the MS COCO dataset [18] to establish two mini datasets, namely *miniCOCO-Alpha* and *miniCOCO-Beta*. Each of them contains 8K images for training and 5K images for validation. We preserve 10 category

labels on *miniCOCO-Alpha*: {*car, handbag, truck, light, bench, chair, horse, bicycle, cup, plant*}, and we preserve an additional 10 category labels on *miniCOCO-Beta*: {*person, bottle, motorcycle, bird, boat, umbrella, sheep, wine glass, table, tv*}. For simplicity, *miniCOCO-Alpha* is denoted as $(\mathbf{I}^\alpha, \mathbf{C}_1^\alpha)$ and *miniCOCO-Beta* is denoted as $(\mathbf{I}^\beta, \mathbf{C}_2^\beta)$, where $\mathbf{I}$, $\mathbf{C}_1$ and $\mathbf{C}_2$ denote the image set, the first 10-category annotation set, and the second 10-category annotation set, respectively. Numbers 1 and 2 correspond to the first 10-category and the second 10-category, respectively. This setting models the central issue when multiple datasets with different category sets are merged.

As we discussed in Section 1, the key step in cross-dataset training for object detection involves seeking the accurate annotation sets $\mathbf{C}_2^\alpha$ and $\mathbf{C}_1^\beta$, and thus, transforming the cross-dataset $(\{\mathbf{I}^\alpha, \mathbf{I}^\beta\}, \{\mathbf{C}_1^\alpha, \mathbf{C}_2^\beta\})$ into a complete form $(\{\mathbf{I}^\alpha, \mathbf{I}^\beta\}, \{\mathbf{C}_1^\alpha, \mathbf{C}_2^\alpha, \mathbf{C}_1^\beta, \mathbf{C}_2^\beta\})$. Next, we conduct a detailed experimental analysis on the generation of the missing annotation sets, $\mathbf{C}_2^\alpha$ and $\mathbf{C}_1^\beta$, with high quality.

**Training Details.** ResNet-50 [32] pretrained on ImageNet [33] is chosen as the backbone of the SSD in this section. During training, we use five image scales {448, 480, 512, 544, 576} (the aspect ratio of the image is 1:1) randomly. The network is trained using the stochastic gradient descent (SGD) algorithm for 100 epochs with 0.9 momentum, 0.0005 weight decay, and 32 batch sizes on two NVIDIA V100 GPUs. The initial learning rate is 0.0026 and it decays at epochs 66 and 83. The loss functions used during training are the cross-entropy loss for the classification branch and the smooth-L1 loss for the regression branch.

**Inference Details.** During the inference phase, the input image is resized to 576 × 576, after which it is input into the entire network to output the predicted bounding boxes with a predicted category. To filter out the redundant background bounding boxes, the confidence threshold is set to 0.01, and non-maximum suppression (NMS) is applied, with an IoU threshold of 0.6 per class, when evaluating the detection performance of the network.

### *3.2. Statistical Characteristic*

A submodel trained using *miniCOCO-Alpha* generates the first 10-category annotation set for *miniCOCO-Beta*, which is denoted as $\mathbf{iC}_1^\beta$. Correspondingly, another submodel trained using *miniCOCO-Beta* generates the second 10-category annotation set for *miniCOCO-Alpha*, which is denoted as $\mathbf{iC}_2^\alpha$. It is worth noting that each element in $\mathbf{iC}_2^\alpha$ and $\mathbf{iC}_1^\beta$ consists of bounding box coordinates, a category label, and the corresponding confidence score. Additionally, the confidence score of each element should be larger than the threshold $T_c$. The next conventional operation involves integrating $\mathbf{iC}_2^\alpha$ and $\mathbf{iC}_1^\beta$ into the datasets and establishing a complete one:

($\{\mathbf{I}^\alpha, \mathbf{I}^\beta\}, \{\mathbf{C}_1^\alpha, \mathbf{iC}_2^\alpha, \mathbf{iC}_1^\beta, \mathbf{C}_2^\beta\}$). Soft-label training is the usual choice for this type of "noisy-annotated" dataset. During training, the loss function for the classification branch is expressed as follows: $L_{cls} = -\sum_{c=0}^{20} y_c \log p_c$. The label for a pseudo-annotation of category $j$ and confidence $s$ is a soft form $\boldsymbol{y} = [y_0, y_1, ..., y_{20}]^T$, where $y_0 = 1 - s$ and $y_j = s$. Through the dataset ($\{\mathbf{I}^\alpha, \mathbf{I}^\beta\}, \{\mathbf{C}_1^\alpha, \mathbf{iC}_2^\alpha, \mathbf{iC}_1^\beta, \mathbf{C}_2^\beta\}$), the soft-label-training submodel achieves 34.4% mAP on the validation set (5K images from *miniCOCO-Alpha* and 5K images from *miniCOCO-Beta*).

Because *miniCOCO-Alpha* and *miniCOCO-Beta* are both sampled from MS COCO, we can obtain the true annotations directly, namely $\mathbf{tC}_2^\alpha$ and $\mathbf{tC}_1^\beta$, where $\mathbf{t}$ denotes true. Using the dataset ($\{\mathbf{I}^\alpha, \mathbf{I}^\beta\}, \{\mathbf{C}_1^\alpha, \mathbf{tC}_2^\alpha, \mathbf{tC}_1^\beta, \mathbf{C}_2^\beta\}$), a theoretically optimal SSD model is trained, and it achieves 38.2% mAP on the same validation set.

We then analyze the predictions of both models on the validation set. As shown in Figure 2, the statistical characteristics have attracted our attention. The quantity of true positives (TP) from the soft-label-training model is similar to that from the theoretically optimal model (the quantity from the theoretically optimal model is approximately 5% higher than that from the soft-label-training model), which means that the maximum recall between the two models is close. However, when we observe the quantities of false positives (FP), we find that the theoretically optimal model outputs much more FP than the soft-label-training model, which is counter-intuitive.

We suggest that the performance gap between the two models is caused by the distribution of detection predictions rather than the quantity of detection outputs. Figure 3 shows the output distribution of category *person*, where the figure on the left is for the theoretically optimal model and the one on the right is for the soft-label-training model. The TP from the theoretically optimal model tend to have a higher confidence score than those of the soft-label-training model. In contrast, more TP from the soft-label-training model are concentrated in the low-score region and severely mixed up with FP, which should be the direct cause of the performance gap between the two detectors. Quantitatively, the TP of the theoretically optimal model that have a confidence score higher than 0.2 account for 68.2% of all TP, whereas those of the soft-label-training model account for a mere 36.3%. As for FP with a confidence score higher than 0.2, their number in the theoretically optimal model is almost four times that in the soft-label-training model.

To verify the difference between the two models, we conduct hard-label training on the dataset ($\{\mathbf{I}^\alpha, \mathbf{I}^\beta\}, \{\mathbf{C}_1^\alpha, \mathbf{iC}_2^\alpha, \mathbf{iC}_1^\beta, \mathbf{C}_2^\beta\}$) for further analysis. The true positive distributions of category *person* for the theoretically optimal model, soft-label-training model, and hard-label-training model are illustrated in Figure 4. The distribution curves of the theoretically optimal model and the hard-label-training model are sim-

ilar. They output more TP in the high-score region, and the peaks of their curves are more to the right than the peak of the soft-label-training model's curve. The TP of the hard-label-training model that have a confidence score higher than 0.2 account for 64.7% of all TP, which is close to that of the theoretically optimal model (68.2%). Additionally, the high-scoring ($>0.9$) TP account for 6.8% and 9.5% for the theoretically optimal model and the hard-label-training model, respectively, which is higher than that of the soft-label-training model (2.1%). Based on these statistics, we suggest that the prediction of the hard-label-training model (the theoretically optimal model is also trained using hard labels) has a higher mean score and contains more FP, whereas that of the soft-label-training model has a lower mean score and contains fewer FP.

### 3.3. Update on Supervisor

Considering the different characteristics of hard-label training and soft-label training, we propose taking advantage of these two strategies to improve the quality of annotations $\mathbf{iC}_2^\alpha$ and $\mathbf{iC}_1^\beta$ from two perspectives. 1) We can *expand* two annotation sets by adding new predictions that can increase the recall of annotations. 2) We can *shrink* two annotation sets by filtering out unreliable predictions to improve the precision of annotations.

Specifically, images from $\mathbf{I}^\alpha$ and $\mathbf{I}^\beta$ are input into the hard-label-training model and the soft-label-training model to obtain new detection results, $\left(\mathbf{H}_2^\alpha, \mathbf{H}_1^\beta\right)$ and $\left(\mathbf{S}_2^\alpha, \mathbf{S}_1^\beta\right)$, respectively. $\mathbf{H}$ denotes the predictions produced using the hard-label-training model, and $\mathbf{S}$ denotes the predictions produced using the soft-label-training model. Therefore the expanding operation using hard-label training is formulated as follows:

$$\begin{cases} \mathbf{hC}_2^{\alpha+} = \mathbf{iC}_2^\alpha \cup (\mathbf{H}_2^\alpha \setminus (\mathbf{C}_1^\alpha \cup \mathbf{iC}_2^\alpha)) \\ \mathbf{hC}_1^{\beta+} = \mathbf{iC}_1^\beta \cup (\mathbf{H}_1^\beta \setminus (\mathbf{C}_2^\beta \cup \mathbf{iC}_1^\beta)) \end{cases} \tag{1}$$

where $\mathbf{hC}_2^{\alpha+}$ denotes the expanded annotation set on dataset *miniCOCO-Alpha* and $\mathbf{hC}_1^{\beta+}$ denotes that on *miniCOCO-Beta*. The operation $\setminus$ minuses the annotation in $\mathbf{H}_2^\alpha$, which has an IoU larger than a threshold $T_e$, with any annotation in $(\mathbf{C}_1^\alpha \cup \mathbf{iC}_2^\alpha)$. Similarly, the expanding operation using soft-label training is formulated as follows:

$$\begin{cases} \mathbf{sC}_2^{\alpha+} = \mathbf{iC}_2^\alpha \cup (\mathbf{S}_2^\alpha \setminus (\mathbf{C}_1^\alpha \cup \mathbf{iC}_2^\alpha)) \\ \mathbf{sC}_1^{\beta+} = \mathbf{iC}_1^\beta \cup (\mathbf{S}_1^\beta \setminus (\mathbf{C}_2^\beta \cup \mathbf{iC}_1^\beta)) \end{cases} \tag{2}$$

| Annotations | $\Delta$Recall(%) | $\Delta$Precision(%) | mAP(%) |
|---|---|---|---|
| $\mathbf{iC}_2^{\alpha}$, $\mathbf{iC}_1^{\beta}$ | 0 | 0 | 34.4 |
| $\mathbf{hC}_2^{\alpha+}$, $\mathbf{hC}_1^{\beta+}$ | +18.8 | -11.8 | **34.8** |
| $\mathbf{sC}_2^{\alpha+}$, $\mathbf{sC}_1^{\beta+}$ | +9.1 | -0.7 | 34.5 |
| $\mathbf{hC}_2^{\alpha-}$, $\mathbf{hC}_1^{\beta-}$ | -5.1 | +12.8 | 33.7 |
| $\mathbf{sC}_2^{\alpha-}$, $\mathbf{sC}_1^{\beta-}$ | -11.6 | **+27.6** | 34.4 |

Table 1: The relative variation of recall and precision when models expand or shrink the initial pseudo annotation sets ($T_e = 0.6$, $T_s = 0.6$). Performance of detection models trained on different pseudo annotation sets. "+" and "−" denote expanding and shrinking operations, respectively. "**h**" and "**s**" denote hard-label training and soft-label training, respectively.

We formulate the shrinking operation using hard-label training as follows:

$$\begin{cases} \mathbf{hC}_2^{\alpha-} = \mathbf{iC}_2^{\alpha} \cap (\mathbf{H}_2^{\alpha} \setminus \mathbf{C}_1^{\alpha}) \\ \mathbf{hC}_1^{\beta-} = \mathbf{iC}_1^{\beta} \cap (\mathbf{H}_1^{\beta} \setminus \mathbf{C}_2^{\beta}) \end{cases} \tag{3}$$

where $\mathbf{hC}_2^{\alpha-}$ denotes the shrunk annotation set on dataset *miniCOCO-Alpha* and $\mathbf{hC}_1^{\beta-}$ denotes that on *miniCOCO-Beta*. The operation $\cap$ preserves the annotation in $\mathbf{iC}_2^{\alpha}$, which has an IoU larger than a threshold $T_s$, with any annotation in $(\mathbf{H}_2^{\alpha} \setminus \mathbf{C}_1^{\alpha})$. It is inspired by the intuition that the bounding boxes detected using both detectors are more reliable. Similarly, the shrinking operation using soft-label training is formulated as follows:

$$\begin{cases} \mathbf{sC}_2^{\alpha-} = \mathbf{iC}_2^{\alpha} \cap (\mathbf{S}_2^{\alpha} \setminus \mathbf{C}_1^{\alpha}) \\ \mathbf{sC}_1^{\beta-} = \mathbf{iC}_1^{\beta} \cap (\mathbf{S}_1^{\beta} \setminus \mathbf{C}_2^{\beta}) \end{cases} \tag{4}$$

Based on the true annotations $\mathbf{tC}_2^{\alpha}$ and $\mathbf{tC}_1^{\beta}$, we can obtain the relative variation of recall and precision for new annotation sets when we expand or shrink them using different models. The results are listed in Table 1, and they confirm our conclusion that the hard-label-training model is conducive to improving the recall of annotation sets, whereas the soft-label-training model is useful for improving the precision of annotation sets. The performance of the models trained using different annotation sets is also shown in Table 1.

## 4. Dynamic Supervisor

According to the analysis conducted in Section 3, a single operation (expanding or shrinking) on the annotation set cannot improve its recall and precision simulta-

neously. Inspired by the characteristics of hard-label training and soft-label training, we delve into a method for improving the quality of the annotation set, and we propose a dynamic supervisor framework to produce a more complete annotation set progressively.

### 4.1. The Overall Framework

The structure of the proposed dynamic supervisor framework is shown in Figure 5, and two datasets annotated with different category sets are merged in our illustration (the notations in Section 3 are followed). There are three steps for generating the final annotation set. The first step of the dynamic supervisor framework is similar to the one employed in previous studies [16, 8], where two detection models are trained using *miniCOCO-Alpha* or *miniCOCO-Beta* individually, after which each detection model generates an initial annotation set for images from the other dataset (*Initial labeling*, as shown in Figure 5).

By combining the ground truths ($\mathbf{C}_1^\alpha$, $\mathbf{C}_2^\beta$) with the generated annotation sets ($\mathbf{iC}_2^\alpha$, $\mathbf{iC}_1^\beta$), *miniCOCO-Alpha* and *miniCOCO-Beta* can be merged. One optional structure of the dynamic supervisor framework is shown in Figure 5 (a), where two datasets are merged after the generation of the initial annotation sets. Next, in the second step (*Expanding* in Figure 5), a detection model trained using ($\{\mathbf{I}^\alpha, \mathbf{I}^\beta\}, \{\mathbf{C}_1^\alpha, \mathbf{iC}_2^\alpha, \mathbf{iC}_1^\beta, \mathbf{C}_2^\beta\}$) generates new annotations through hard-label training to expand the initial annotation sets ($\mathbf{iC}_2^\alpha$, $\mathbf{iC}_1^\beta$) into ($\mathbf{hC}_2^{\alpha+}$, $\mathbf{hC}_1^{\beta+}$). In the third step (*Shrinking* in Figure 5), another detection model trained using ($\{\mathbf{I}^\alpha, \mathbf{I}^\beta\}, \{\mathbf{C}_1^\alpha, \mathbf{hC}_2^{\alpha+}, \mathbf{hC}_1^{\beta+}, \mathbf{C}_2^\beta\}$) filters out unreliable annotations through soft-label training to shrink the expanded annotation sets ($\mathbf{hC}_2^{\alpha+}$, $\mathbf{hC}_1^{\beta+}$) into ($\mathbf{sC}_2^{\alpha-}$, $\mathbf{sC}_1^{\beta-}$). Based on the dataset ($\{\mathbf{I}^\alpha, \mathbf{I}^\beta\}, \{\mathbf{C}_1^\alpha, \mathbf{sC}_2^{\alpha-}, \mathbf{sC}_1^{\beta-}, \mathbf{C}_2^\beta\}$), we can obtain the final model for cross-dataset object detection.

In the first structure of the dynamic supervisor discussed above, the detection models (in the second and third steps) must update their supervision information independently. Considering the risk that the models are prone to overfitting to noisy annotations in this "self-annotated mechanism", we propose another structure, which is a "cross-annotated mechanism," as shown in Figure 5 (b). Unlike in the first structure, the two datasets are not merged after the generation of the initial annotation sets in the first step. When the two original datasets are augmented with the generated annotation sets ($\mathbf{I}^\alpha, \mathbf{C}_1^\alpha, \mathbf{iC}_2^\alpha$) and ($\mathbf{I}^\beta, \mathbf{iC}_1^\beta, \mathbf{C}_2^\beta$), they are responsible for training two models individually through hard-label training. The two hard-label-training models then generate new annotations for images from the other dataset to expand their initial annotation set in the second step. Because these two models are trained using the augmented datasets, their detection distributions are different from those of the previous two models in the first step. Therefore, they can detect objects

| Setting | Datasets |
|:---:|:---:|
| A | miniCOCO-Alpha [18] + miniCOCO-Beta [18] |
| B | VOC [17] + COCO [18] (w/o VOC categories) |
| C | VOC [17] + SUN-RGBD [19] + LISA-Signs [20] |

Table 2: Three combinations of datasets in our experiments. Setting A is the dataset combination in Section 3. B and C are the same setting as that of [8].

that are ignored, thereby improving the recall of the annotation sets. Although the recall is improved, the expanding operation probably results in a decrease in precision. Accordingly, two other models are then trained through soft-label training using the expanded annotation sets $((\mathbf{I}^{\alpha}, \mathbf{C}_1^{\alpha}, \mathbf{hC}_2^{\alpha+})$ or $(\mathbf{I}^{\beta}, \mathbf{hC}_1^{\beta+}, \mathbf{C}_2^{\beta}))$. The new annotations generated using two soft-label-training models are used to filter out the unreliable annotations of the expanded annotation sets in the third step. The objects detected through the two different models are more likely to be TP than FP. Therefore, there will be an improvement in the precision of annotation sets after the shrinking operation.

Finally, *miniCOCO-Alpha* and *miniCOCO-Beta* are merged, and the final model is trained using this multiple-updated dataset: $(\{\mathbf{I}^{\alpha}, \mathbf{I}^{\beta}\}, \{\mathbf{C}_1^{\alpha}, \mathbf{sC}_2^{\alpha-}, \mathbf{sC}_1^{\beta-}, \mathbf{C}_2^{\beta}\})$. We suppose that the recall and precision of the annotation sets can be promoted comprehensively through the multiple update steps.

### 4.2. Experimental Verification

**Cross-dataset Setting.** To verify the effectiveness of the proposed dynamic supervisor framework, we define three different combinations of datasets that cover diverse scenarios. In Table 2, we summarize the dataset settings used in our experiments. Setting A contains two mini datasets sampled from MS COCO, which we used in Section 3. Setting B combines PASCAL VOC [17] (20 categories) and MS COCO [18] (without VOC categories) to verify the performance of the proposed method on large-scale datasets. Setting C consists of three datasets from different scenarios: PASCAL VOC is a general dataset containing 20 common categories, SUN-RGBD [19] includes 18 categories of indoor scenes, and LISA-Signs [20] is a driving scenes dataset, which contains 4 traffic signs. This setting combines multiple datasets with large gaps, and frequent scene overlapping exists among these datasets.

**Comparison with Other Methods.** We apply the proposed dynamic supervisor framework to the Faster-RCNN [36] model. ResNet-50 [32] with FPN [37] is used as the backbone, and it is pretrained on ImageNet [33]. We follow the implementation presented by [38], where input images are resized to keep their shorter

| Method | COCO | MIX |
|---|---|---|
| Naive-combination | 42.6 | 43.7 |
| Partial-loss [34] | 43.6 | 44.6 |
| UOD [35] + Merge | 45.6 | 46.1 |
| Pseudo-Labeling [8] | 50.3 | 52.2 |
| Static Supervisor (ours) | 53.3 | 51.5 |
| Dynamic Supervisor (ours) | **56.2** | **55.8** |

Table 3: Detection performance(mAP in %) of setting B on different validation sets. The backbone of all methods is ResNet-50.

| Method | backbone | mAP |
|---|---|---|
| Naive-combination | ResNet-50 | 58.1 |
| Partial-loss [34] | ResNet-50 | 58.3 |
| UOD [35] + Merge | ResNet-50 | 59.3 |
| Min-Entropy loss [8] | ResNet-50 | 58.7 |
| Pseudo-Labeling [8] | ResNet-50 | 61.1 |
| Static Supervisor (ours) | ResNet-50 | 60.7 |
| Dynamic Supervisor (ours) | ResNet-50 | **61.7** |

Table 4: Detection performance(mAP in %) on setting C.

side as one of $\{640, 672, 704, 736, 768, 800\}$ randomly and their longer side less than or equal to 1,333 during training.

For setting B, we evaluate our model on two datasets: COCO and MIX [8] (a combination of VOC and COCO). The experimental results are shown in Table 3. Compared with the naive combination of the two datasets (the first row in Table 3), the partial loss method proposed by [34] only takes a small step forward. The following pseudo-labeling method proposed by [8] achieves much better performance (50.3% *vs.* 42.6%, 52.2% *vs.* 43.7%), which demonstrates the effectiveness of pseudo-annotation. The performance of the static supervisor is comparable to that of the one proposed by [8] (53.3% *vs.* 50.3%, 51.5% *vs.* 52.2%). However, the dynamic supervisor framework goes further and achieves a new state-of-the-art performance (56.2% on COCO and 55.8% on MIX). This performance demonstrates the effectiveness of the dynamic supervisor framework in a large-scale dataset setting.

For setting C, there are four overlapping categories between PASCAL VOC and SUN-RGBD. Therefore, the entire merged dataset has 38 categories in total. A vali-

dation set [8] of 1,500 images (from three datasets) is annotated for all 38 categories for evaluation. The results of the experiments conducted in this setting are listed in Table 4. The static supervisor framework achieves a performance that is comparable to that of the pseudo-labeling method [8] (60.7% *vs.* 61.1%). Based on the static supervisor framework, the proposed dynamic supervisor framework enhances the performance further by 1% in mAP (61.7% *vs.* 60.7%). The best performance in this setting demonstrates the generality of our method in a complicated setting (multiple datasets with frequent scene overlapping).

Comparing the promotions brought by the dynamic supervisor in setting B and setting C, the performance superiority in setting B is larger than that in setting C. This inapparent superiority in setting C results from the slight annotation missing. The more annotations are lost, the more promotions the dynamic supervisor will bring. Quantitively, there are 4.8 predicted instances per image for setting B, whereas there are only 1.0 predicted instances per image for setting C. In other words, the training task in setting C is closer to a fully-supervised problem, resulting in minor performance gaps among all competitors.

### 4.3. Ablation Studies

We follow the controlled experimental setting in Section 3.1 to analyze the proposed dynamic supervisor framework. The implementation details are already described in Section 3.1. Ablation results are shown in Table 5 - Table 7, Figure 6, and Figure 7. Details about ablation studies are discussed in the following.

**Dynamic Hyperparameters:** Figure 6 shows the detection performance of models trained using initial annotation sets with a different confidence threshold $T_c$. It is hard for us to determine which threshold is the best (it should be between 0.15 and 0.3 in this setting), especially when encountering other datasets or detection frameworks. Meanwhile, considering the recognition abilities of a detection model vary from one category to another, we propose selecting the confidence threshold adaptively. The confidence threshold for each category is determined when the F1-score of this category reaches the maximum on the validation set. This method achieves the best performance, as shown in Figure 6, which is simple but effective.

Figure 7 shows the detection performance for different choices of $T_e$ and $T_s$. All these hyperparameter combinations can improve the detection performance compared with the baseline model (34.4% mAP). From Eq. (1) and (4), we can observe that a larger $T_e$ introduces more new pseudo-annotations in the expanding operation. Additionally, a smaller $T_s$ preserves more old pseudo-annotations in the shrinking operation. As shown in Figure 7, it is better to apply a large $T_e$ to obtain more recall improvement and a small $T_s$ to avoid many true annotations being filtered out. In

| Method | Pseudo-label | Expand | Shrink | mAP |
|---|---|---|---|---|
| Naive | | | | 26.6 |
| Static | ✓ | | | 34.4 |
| Static+expand | ✓ | ✓ | | 35.0 |
| Static+shrink | ✓ | | ✓ | 34.8 |
| Dynamic | ✓ | ✓ | ✓ | **35.5** |

Table 5: Detection performance (mAP in %) of models with different operations. *"Pseudo-label"* means using submodels to generate missing annotations. *"Expand"* and *"Shrink"* mean updating on the annotation sets.

| Operation | $\Delta$Recall(%) | $\Delta$Precision(%) | mAP |
|---|---|---|---|
| Pseudo-label | 0 | 0 | 34.4 |
| Expand | +27.6 | -5.2 | 35.0 |
| Expand-shrink | +16.8 | +25.9 | **35.5** |
| Shrink | -6.0 | +36.6 | 34.8 |
| Shrink-expand | +23.1 | +9.8 | 35.2 |

Table 6: The relative variation of recall and precision when new submodels expand or shrink the original pseudo annotations set.

the remainder of this paper, we use $T_e = 0.7$ and $T_s = 0.5$ for ablation studies.

**Dynamic Supervisor *vs.* Static Supervisor:** We report the detection performance of models with different operations in Table 5. The model trained using the naive combination of two datasets, without any other strategies, performs the worst. When the static supervisor framework is applied, the corresponding model enhances the performance by 7.8% in mAP, showing the effect of pseudo-labeling on the cross-dataset object detection task. Furthermore, a dynamic supervisor framework is proposed to update the initial annotation sets multiple times using two types of submodels. From the last three models shown in Table 5, the promotions resulting from the expanding operation, the shrinking operation, or both are clear. In Table 6, we record the variation in recall and precision when new submodels apply expanding, shrinking, or a combination of both to the annotation set. A single expanding operation on the annotation sets increases the recall but decreases the precision, and the shrinking operation is simply the opposite. Nevertheless, when we combine the advantages of these two operations and apply them sequentially, both the recall and precision can be improved.

| Type | Expand | Shrink | mAP |
|------|:------:|:------:|-----|
| Self-annotated | ✓ | | 34.8 |
| | ✓ | ✓ | 34.6 |
| Cross-annotated | ✓ | | 35.0 |
| | ✓ | ✓ | **35.5** |

Table 7: Detection performance (mAP in %) of different types of dynamic supervisor. *"Expand"* and *"Shrink"* mean the update operation on pseudo annotations.

**Cross-annotated *vs.* Self-annotated:** In Section 4.1, we introduce two types of the dynamic supervisor framework, which is a "self-annotated mechanism" for the first one and a "cross-annotated mechanism" for the second one. In the self-annotated mechanism, two datasets are merged, and the detection model is trained using this merged dataset to expand or shrink the initial annotation sets after they are generated. In the cross-annotated mechanism, there are two models trained using their respective datasets, which are already augmented with annotation sets, after which they are utilized to expand or shrink the annotation set of each other. The detection performance of these two types of dynamic supervisors is shown in Table 7. When the expanding operation is applied to the annotation set, both mechanisms result in performance improvements. However, when the shrinking operation is applied subsequently, the detection performance of the self-annotated mechanism tends to degrade (mAP decreases from 34.8% to 34.6%). We suggest that the model trained using a noisy dataset is prone to overfitting to incorrect annotations. Consequently, in the self-annotated mechanism, the knowledge learned from such a noisy dataset is that it is difficult to eliminate the noise of this dataset continuously. In contrast, the cross-annotated mechanism avoids this dilemma and can progressively improve performance (mAP increases from 34.4% to 35.5% step by step).

**Operation Sequence:** There are two operations in the proposed dynamic supervisor framework for improving the recall or precision of the annotation set. Here, we attempt to update the annotation set through a different sequence of operations to explore the relationship between quality variation and performance improvement. The results of the experiments on the two types of operation sequences are shown in Table 6. First, and most importantly, the two types of operation sequences can improve the quality of the annotation set and achieve similar detection performance (35.5% *vs.* 35.2%). However, the different operation sequences would result in different improvements in recall and precision. When we first expand the annotation set and then shrink it, we obtain more improvement in precision and less improvement

in recall. This is because the later shrinking operation is prone to missing TP. Therefore, it suppresses recall. Conversely, the improvement in the recall will be higher. The results show the flexibility of the dynamic supervisor framework, and the goal is to improve the quality of the annotation set comprehensively.

### 4.4. Qualitative Visualizations:

Finally, we visualize the pseudo-labeling results in different steps of our dynamic supervisor framework, as shown in Figure 8. These images are sampled from *miniCOCO-Beta*, and their original ground-truth boxes are not labeled here. As described in previous sections, the pseudo-labeling results of a single model are biased, and it is difficult for a single model to achieve high recall and high precision (the first row of Figure 8). Therefore, we propose updating the annotations set multiple times in a dynamic framework. The expanding operation is effective for increasing the number of TP. However, it also facilitates the introduction of new FP (the second row of Figure 8). Consequently, the shrinking operation is applied to discard them and obtain a cleaner annotation set (the third row of Figure 8). However, compared with the ground truth boxes, objects, such as the cup (on the table) and handbag (hanging on the sofa), which are small or hard to distinguish, still need to be found.

## 5. Conclusion

In this study, to address the problem of cross-dataset object detection, we reveal the implicit connections between hard- and soft-label training and the methods for pseudo-annotation ensembling. We show that hard-label training and soft-label training are conducive to improving the recall and precision of the detection results, respectively. Based on this, we propose a dynamic supervisor framework, which polishes the annotations dynamically and selects predictions adaptively based on category. The proposed dynamic supervisor framework updates the annotation set multiple times and improves its quality progressively. Experiments conducted on different combinations of several datasets demonstrate the effectiveness of the proposed dynamic supervisor framework.

## 6. Acknowledgements

## References

[1] M. I. Jordan, D. E. Rumelhart, Forward models: Supervised learning with a distal teacher, Cognitive science 16 (3) (1992) 307–354. 1

[2] L. Bossard, M. Guillaumin, L. Van Gool, Food-101–mining discriminative components with random forests, in: European conference on computer vision, Springer, 2014, pp. 446–461. 1

[3] L. Fei-Fei, R. Fergus, P. Perona, Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories, in: 2004 conference on computer vision and pattern recognition workshop, IEEE, 2004, pp. 178–178. 1

[4] W. Li, R. Zhao, T. Xiao, X. Wang, Deepreid: Deep filter pairing neural network for person re-identification, in: CVPR, 2014. 1

[5] L. Cao, Z. Liu, T. S. Huang, Cross-dataset action detection, in: 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, IEEE, 2010, pp. 1998–2005. 2

[6] P. Peng, T. Xiang, Y. Wang, M. Pontil, S. Gong, T. Huang, Y. Tian, Unsupervised cross-dataset transfer learning for person re-identification, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 1306–1315. 2

[7] T. Perrett, D. Damen, Recurrent assistance: cross-dataset training of lstms on kitchen tasks, in: Proceedings of the IEEE International Conference on Computer Vision Workshops, 2017, pp. 1354–1362. 2

[8] X. Zhao, S. Schulter, G. Sharma, Y.-H. Tsai, M. Chandraker, Y. Wu, Object detection with a unified label space from multiple datasets, in: European Conference on Computer Vision (ECCV), 2020. 2, 4, 9, 10, 11, 12

[9] M. Oquab, L. Bottou, I. Laptev, J. Sivic, Is object localization for free?-weakly-supervised learning with convolutional neural networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 685–694. 2

[10] P. Tang, X. Wang, S. Bai, W. Shen, X. Bai, W. Liu, A. Yuille, Pcl: Proposal cluster learning for weakly supervised object detection, IEEE transactions on pattern analysis and machine intelligence 42 (1) (2018) 176–191. 2

[11] W. Jiang, Z. Zhao, F. Su, Y. Fang, Dynamic proposal sampling for weakly supervised object detection, Neurocomputing (2021). 2

[12] S. Yi, H. Ma, X. Li, Y. Wang, Wsodpb: Weakly supervised object detection with pcsnet and box regression module, Neurocomputing 418 (2020) 232–240. 2

[13] O. Chapelle, B. Scholkopf, A. Zien, Semi-supervised learning (chapelle, o. et al., eds.; 2006)[book reviews], IEEE Transactions on Neural Networks 20 (3) (2009) 542–542. 2

[14] J. Gao, J. Wang, S. Dai, L.-J. Li, R. Nevatia, Note-rcnn: Noise tolerant ensemble rcnn for semi-supervised object detection, in: Proceedings of the IEEE international conference on computer vision, 2019, pp. 9508–9517. 2

[15] X. Zhu, A. B. Goldberg, Introduction to semi-supervised learning, Synthesis lectures on artificial intelligence and machine learning 3 (1) (2009) 1–130. 2

[16] A. Rame, E. Garreau, H. Ben-Younes, C. Ollion, Omnia faster r-cnn: Detection in the wild through dataset merging and soft distillation, arXiv preprint arXiv:1812.02611 (2018). 2, 4, 9

[17] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, A. Zisserman, The pascal visual object classes (voc) challenge, International journal of computer vision 88 (2) (2010) 303–338. 3, 10

[18] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, C. L. Zitnick, Microsoft coco: Common objects in context, in: European conference on computer vision, Springer, 2014, pp. 740–755. 3, 4, 10

[19] S. Song, S. P. Lichtenberg, J. Xiao, Sun rgb-d: A rgb-d scene understanding benchmark suite, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 567–576. 3, 10

[20] A. Mogelmose, M. M. Trivedi, T. B. Moeslund, Vision-based traffic sign detection and analysis for intelligent driver assistance systems: Perspectives and survey, IEEE Transactions on Intelligent Transportation Systems 13 (4) (2012) 1484–1497. 3, 10

[21] I. Radosavovic, P. Dollár, R. Girshick, G. Gkioxari, K. He, Data distillation: Towards omni-supervised learning, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 4119–4128. 3

[22] J. Jeong, S. Lee, J. Kim, N. Kwak, Consistency-based semi-supervised learning for object detection (2019). 3

[23] K. Wang, L. Lin, X. Yan, Z. Chen, D. Zhang, L. Zhang, Cost-effective object detection: Active sample mining with switchable selection criteria, IEEE transactions on neural networks and learning systems 30 (3) (2018) 834–850. 3

[24] Z. Wu, N. Bodla, B. Singh, M. Najibi, R. Chellappa, L. S. Davis, Soft sampling for robust object detection, arXiv preprint arXiv:1806.06986 (2018). 3

[25] Y. Yang, K. J. Liang, L. Carin, Object detection as a positive-unlabeled problem, arXiv preprint arXiv:2002.04672 (2020). 4

[26] S. Chadwick, P. Newman, Training object detectors with noisy data, in: 2019 IEEE Intelligent Vehicles Symposium (IV), IEEE, 2019, pp. 1319–1325. 4

[27] M. Abbasi, D. Laurendeau, C. Gagne, Self-supervised robust object detectors from partially labelled datasets, arXiv preprint arXiv:2005.11549 (2020). 4

[28] Y. Zou, Z. Yu, X. Liu, B. Kumar, J. Wang, Confidence regularized self-training, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2019, pp. 5982–5991. 4

[29] R. Müller, S. Kornblith, G. E. Hinton, When does label smoothing help?, in: NeurIPS, 2019. 4

[30] M. Lukasik, S. Bhojanapalli, A. Menon, S. Kumar, Does label smoothing mitigate label noise?, in: International Conference on Machine Learning, PMLR, 2020, pp. 6448–6458. 4

[31] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, A. C. Berg, Ssd: Single shot multibox detector, in: European conference on computer vision, Springer, 2016, pp. 21–37. 4

[32] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778. 5, 10

[33] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al., Imagenet large scale visual recognition challenge, International journal of computer vision 115 (3) (2015) 211–252. 5, 10

[34] T. Cour, B. Sapp, B. Taskar, Learning from partial labels, The Journal of Machine Learning Research 12 (2011) 1501–1536. 11

[35] X. Wang, Z. Cai, D. Gao, N. Vasconcelos, Towards universal object detection by domain attention, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 7289–7298. 11

[36] S. Ren, K. He, R. Girshick, J. Sun, Faster r-cnn: Towards real-time object detection with region proposal networks, in: Advances in neural information processing systems, 2015, pp. 91–99. 10

[37] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, S. Belongie, Feature pyramid networks for object detection, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 2117–2125. 10

[38] Y. Wu, A. Kirillov, F. Massa, W.-Y. Lo, R. Girshick, Detectron2, https://github.com/facebookresearch/detectron2 (2019). 10
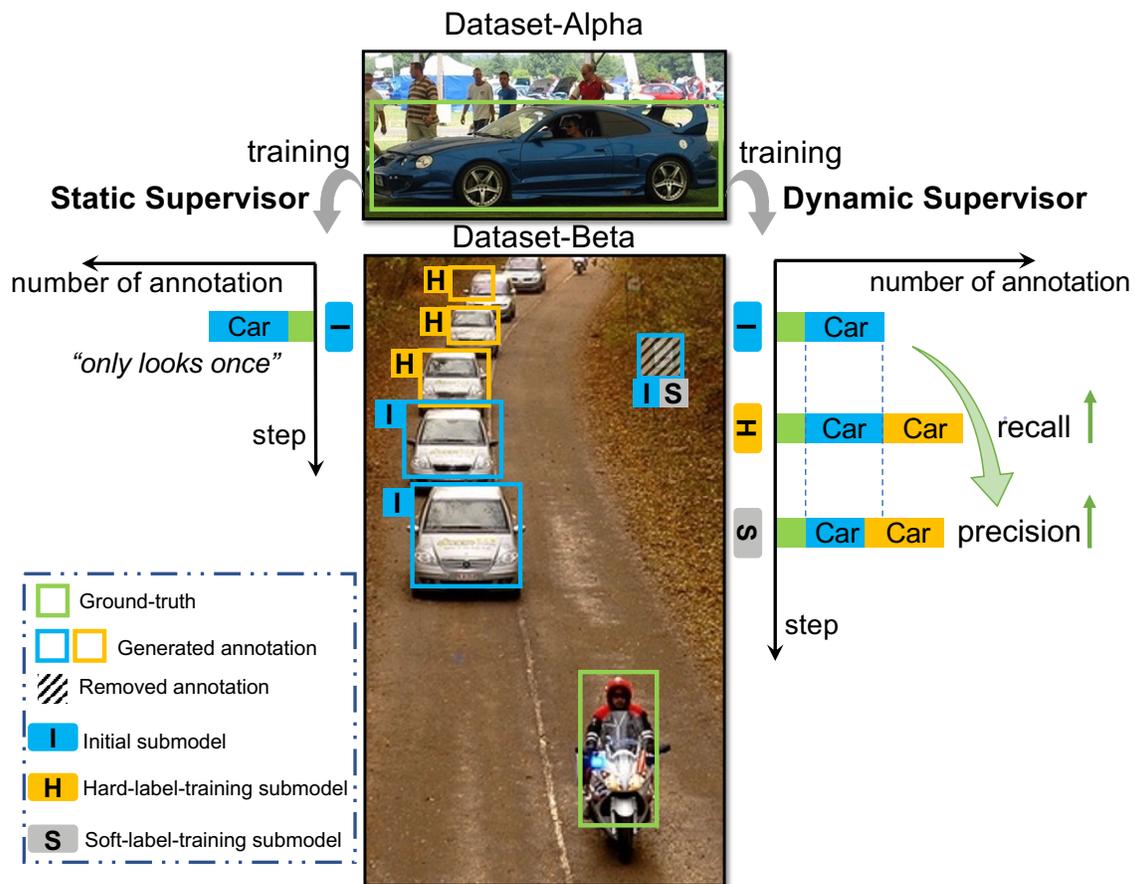
Figure 1: The difference between static supervisor and dynamic supervisor. There is no car category in the annotation of dataset Beta. Submodels trained on dataset Alpha are used to generate the missing annotations of dataset Beta. For the static supervisor, the initial submodel only looks once and detects three cars (one of them is false) as annotations. For the dynamic supervisor, after the annotation generation of the initial submodel, a hard-label-training submodel generates new annotations to increase the recall rate. Then, a soft-label-training submodel is utilized to filter out the unreliable annotations that increases the precision rate.

| category | Optimal | Soft | $P = \dfrac{Optimal - Soft}{Soft}$ (%) |
|---|---|---|---|
| Person | 3.92k | 3.72k | ▇ 5.38% |
| | 265k | 227k | ▇▇▇▇ 16.74% |
| Bottle | 2.49k | 2.38k | ▇ 4.62% |
| | 322k | 245k | ▇▇▇▇▇▇ 31.43% |
| Handbag | 1.38k | 1.30k | ▇ 6.15% |
| | 254k | 260k | ▇ -2.30% |
| Bird | 1.82k | 1.86k | ▇ -2.15% |
| | 407k | 374k | ▇▇ 8.82% |
| All | 28.9k | 27.6k | ▇ 4.71% |
| | 4325k | 3100k | ▇▇▇▇▇▇▇ 39.52% |

▇ TP ▇ FP

Figure 2: TP (in green color) and FP (in red color) quantities of detection outputs for the theoretically optimal model and soft-labeling-training model. The statistics of 4 specific categories are listed here. "All" denotes the total quantities of 20 categories. $P$ reflects that the TP quantities of two models are similar but the FP quantity of the theoretically optimal model is much more than that of the soft-labeling-training model.
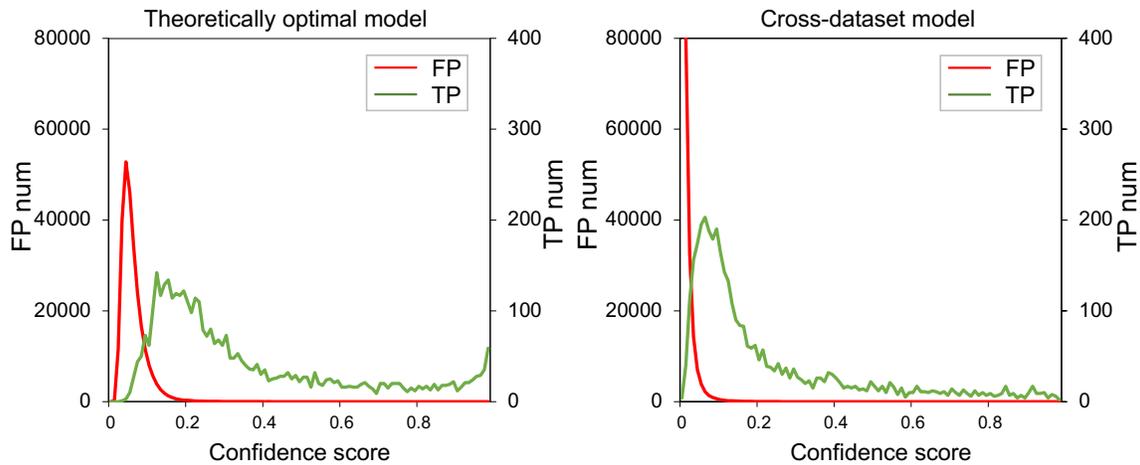


Figure 3: The numbers of detection results with different confidence scores of TP and FP respectively. **Left**: detection output distribution of the theoretically optimal model. **Right**: detection output distribution of the cross-dataset model (soft-label training).
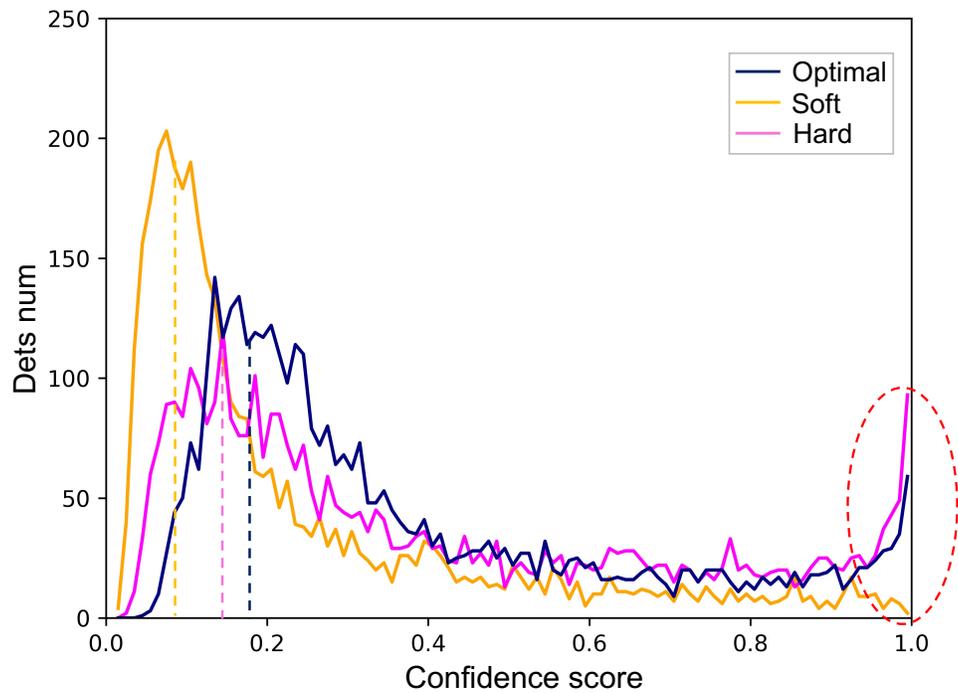
Figure 4: The numbers of TP with different confidence scores for three models. *"Optimal"* is the curve of the theoretically optimal model. *"Soft"* and *"Hard"* is the curve of the soft-label-training model and hard-label-training model respectively.
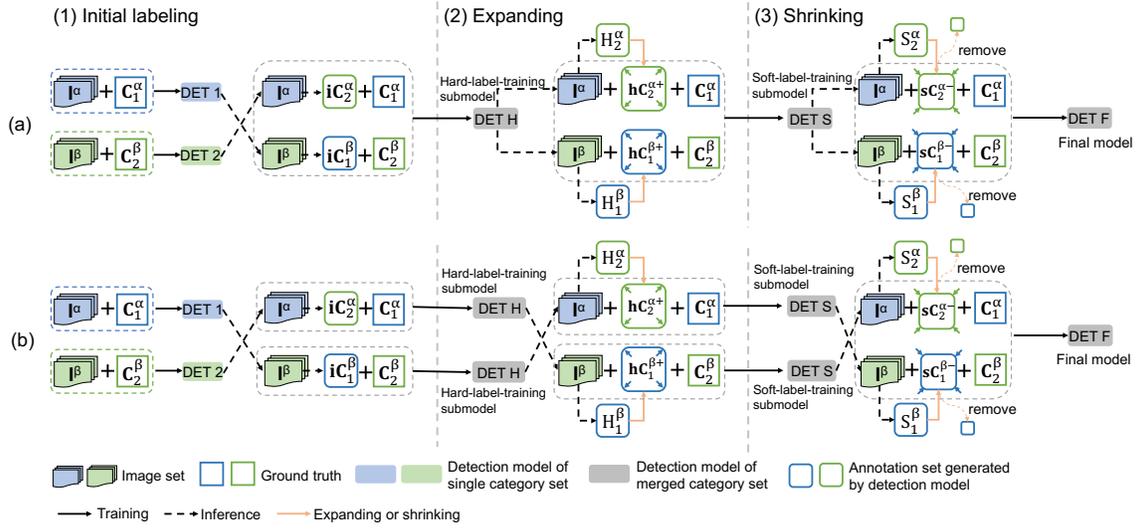
Figure 5: The structure of dynamic supervisor framework. (a) Self-annotated mechanism. (b) Cross-annotated mechanism. The notations in Section 3 are followed here. $\alpha$ and $\beta$ correspond to different image sets and numbers 1 and 2 correspond to different category sets.
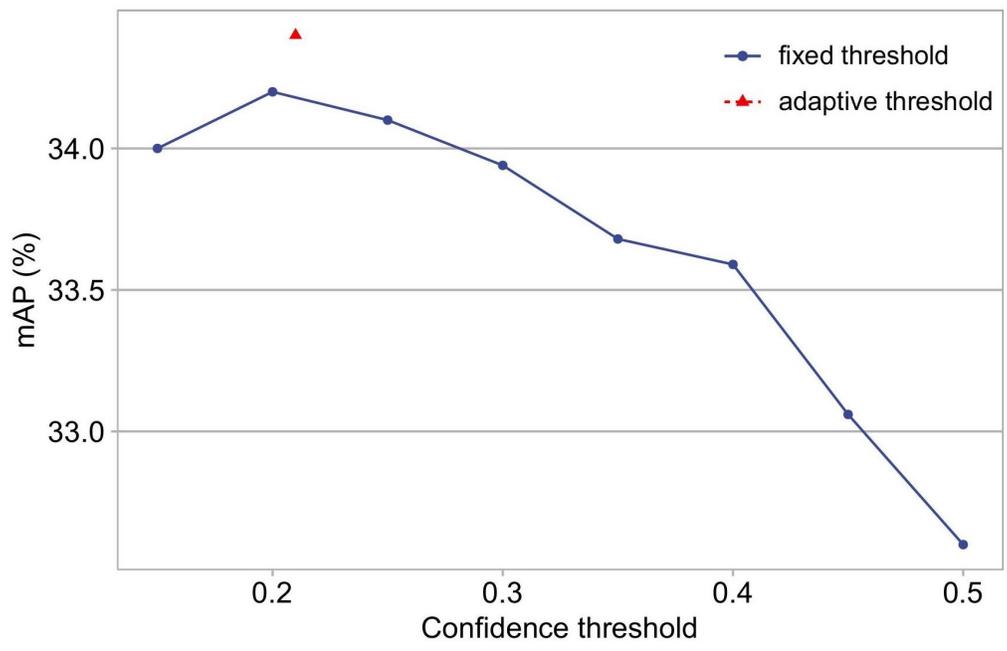
Figure 6: Detection performance (mAP in %) of models trained on initial annotation sets with fixed confidence threshold and adaptive threshold. The average value of adaptive threshold is 0.21 in this setting.
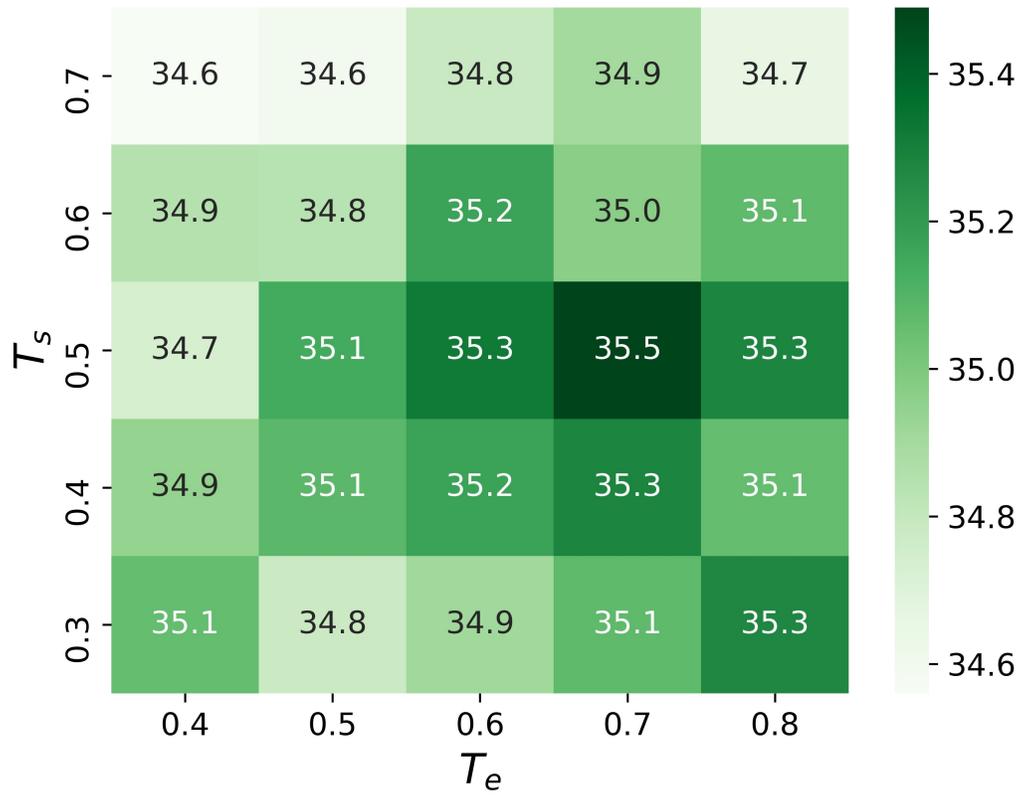
Figure 7: Detection performance (mAP in %) for different choices of $T_e$ and $T_s$.

Figure 8: Visualizations of pseudo-labeling results in different steps of the dynamic supervisor framework. These images are sampled from the *miniCOCO-beta* and their original ground truth boxes are not labeled here. Green rectangle indicates true annotations, red rectangle indicates false annotations, and blue dotted rectangle indicates missing annotations. The numbers of true annotations (TP), false annotations (FP), and missing annotations (FN) are listed inside every image.