

Label Matching Semi-Supervised Object Detection

Binbin Chen², Weijie Chen^{1,2,†}, Shicai Yang², Yunyi Xuan², Jie Song¹
Di Xie², Shiliang Pu², Mingli Song¹, Yueting Zhuang^{1,†}

¹Zhejiang University, ²Hikvision Research Institute

{chenbinbin8, chenweijie5, yangshicai, xuanyunyi, xiedi, pushiliang.hri}@hikvision.com

{sjie, songml, yzhuang}@zju.edu.cn

Abstract

Semi-supervised object detection has made significant progress with the development of mean teacher driven self-training. Despite the promising results, the label mismatch problem is not yet fully explored in the previous works, leading to severe confirmation bias during self-training. In this paper, we delve into this problem and propose a simple yet effective LabelMatch framework from two different yet complementary perspectives, i.e., distribution-level and instance-level. For the former one, it is reasonable to approximate the class distribution of the unlabeled data from that of the labeled data according to Monte Carlo Sampling. Guided by this weakly supervision cue, we introduce a re-distribution mean teacher, which leverages adaptive label-distribution-aware confidence thresholds to generate unbiased pseudo labels to drive student learning. For the latter one, there exists an overlooked label assignment ambiguity problem across teacher-student models. To remedy this issue, we present a novel label assignment mechanism for self-training framework, namely proposal self-assignment, which injects the proposals from student into teacher and generates accurate pseudo labels to match each proposal in the student model accordingly. Experiments on both MS-COCO and PASCAL-VOC datasets demonstrate the considerable superiority of our proposed framework to other state-of-the-arts. Code will be available at <https://github.com/hikvision-research/SSOD>.

1. Introduction

Supervised learning has advanced object detection in the past few years, benefited from tremendous labeled training data [6, 21, 32, 34, 43]. However, it is extremely expensive and time-consuming to collect accurate annotations. As an alternative, semi-supervised object detection (SSOD) is proposed to use a small amount of labeled data in conjunc-

[†]Corresponding author

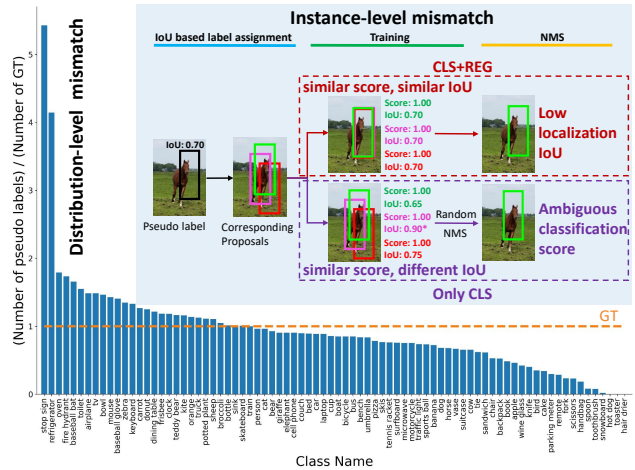


Figure 1. Label mismatch problems on the MS-COCO dataset. 1) **Distribution-level mismatch**: there exists a bias between the pseudo labels produced by the single confidence threshold and the ground truth labels (GT) during self-training, as shown in the relation of the blue bar and the orange dotted line. 2) **Instance-level mismatch**: there are two kinds of training patterns for the unlabeled data in the previous SSOD frameworks. One is the same as supervised learning, using both classification and box regression for optimization, which will overfit the poor-quality pseudo labels and result in low localization accuracy. To avoid incorrect box regression, another one merely exploits a classification objective [27], which will bring ambiguity due to the similar classification scores to confuse the post-processing of Non-Maximum-Suppression (NMS).

tion with a large amount of unlabeled data to optimize the detectors [17, 27, 40, 49, 52]. Recently, SSOD has achieved growing interest in the object detection community.

Self-training has been proven useful in SSOD, especially mean teacher framework [27, 49], which annotates the unlabeled data by a gradually evolving teacher and guides the learning of a student in a mutually beneficial manner. As the key process of mean teacher, the existing pseudo labeling methods [27, 49, 52] simply utilize a hand-crafted confidence threshold to filter out low-quality pseudo labels and

directly treat the remaining ones as reliable pseudo labels. However, it is inevitable to encounter the label mismatch problem, leading to severe confirmation bias [1] during self-training. In this paper, we delve into this problem from two perspectives, *i.e.*, *distribution-level* and *instance-level*.

From the perspective of the distribution-level label mismatch problem, it is extremely difficult to generate unbiased pseudo labels to match the ground-truth labels with consistent class distribution by using a single and fixed confidence threshold due to the class-imbalanced data distribution. As shown in Fig. 1, the number of pseudo labels is much higher than the ground-truth labels in some classes, while far less in some other classes, resulting in abundant false positives and false negatives. From the perspective of the instance-level label mismatch problem, the existing methods directly follow supervised object detection [34] for label assignment. However, the situation is totally different in semi-supervised learning since the quality of pseudo label cannot be guaranteed, leading to label assignment ambiguity problem as illustrated in Fig. 1. Especially in mean teacher driven self-training framework, it is crucial to study how to assign the pseudo labels generated by the mean teacher to the proposals generated by the student network rather than a rough IoU based label assignment manner [34]. Based on the aforementioned challenges of label mismatch in two different yet complementary granularities, we begin our study and develop a *LabelMatch* framework.

To address the first challenge, we present a very simple *re-distribution mean teacher*. Assumed that the labeled data is selected from the entire data gallery via Monte Carlo Sampling. In this way, the label distribution of the unlabeled data can be approximated from that of the labeled data. In fact, we have evaluated the label distribution of the labeled and unlabeled data in several popular SSOD datasets, and they all meet this hypothesis which can be exploited as a weakly supervision cue for pseudo labeling. Under this inspiration, in contrast to a single and fixed confidence threshold, we utilize *adaptive label-distribution-aware confidence thresholds* (ACT) to generate unbiased pseudo labels for the unlabeled data, supervised by the label distribution in the labeled data. The ACT are category-specific and adaptively up-to-date during self-training.

To address the second challenge, we propose a novel *proposal self-assignment* method. Before introducing our method, we should highlight that it is infeasible to set all pseudo labels as hard labels due to the poor quality of the pseudo labels, especially at the beginning of self-training. Under this consideration, we divide the pseudo labels into reliable ones and uncertain ones according to the confidence score. We treat the reliable labels as hard labels for model optimization identically to the supervised manner, while exploiting the uncertain ones via the proposal self-assignment method for soft learning. Detailedly, the proposals from the

student are injected into the teacher for proposals correction, which can provide corresponding soft labels to rectify each proposal accordingly. Besides, to encourage positive feedback during self-training, we introduce a *reliable pseudo label mining* (RPLM) strategy to further improve the performance, which aims to convert the high-quality uncertain pseudo labels into reliable ones in a curriculum way.

We benchmark LabelMatch with the same experimental settings to Unbiased-Teacher [27] using the MS-COCO [25] and PASCAL-VOC [10] datasets, namely *COCO-standard*, *COCO-additional*, and *VOC*. LabelMatch achieves new state-of-the-art results across all benchmarks. Especially in the settings with scarce labeled data, *i.e.*, *COCO-standard* with only 1% labeled data and *VOC*, our method can surpass the previous state-of-the-arts by a large margin.

The contributions of this paper are listed as follows:

- We contribute to analyzing the label mismatch problem from the perspectives of distribution-level and instance-level, which provides a brand-new direction for SSOD.
- We propose a simple yet effective LabelMatch framework to address the label mismatch problems in SSOD. In this framework, we 1) present a re-distribution mean teacher to address the distribution-level label mismatch problem; 2) design a proposal self-assignment scheme to address the instance-level label mismatch problem; 3) introduce a reliable pseudo label mining strategy for pseudo label re-calibration during self-training.
- The LabelMatch framework achieves new state-of-the-arts on many popular SSOD benchmarks. Also, we build a MMDetection-based semi-supervised object detection codebase for the fair study of SSOD algorithms.

2. Related Work

Semi-Supervised Classification. The general methods can be roughly categorized into two types. One is consistency regularization, assuming the model’s predictions to be invariant even if various perturbations are applied. There are different kinds of perturbations, including model-level perturbations [16, 36, 42], image augmentations [45], and adversarial training [28]. Another one is self-training, aka pseudo labeling, which regards the predictions as pseudo labels. For instance, NoisyStudent [46] evolves the pseudo labels for model optimization iteratively. MixMatch [2] uses mixup augmentation and averages different augmented predictions to generate pseudo labels. FixMatch [39] uses the weakly augmented data for pseudo labeling while exploiting the strongly augmented data for model training.

Semi-Supervised Object Detection. The technologies in SSOD are inherited from semi-supervised classification, dividing into consistency regularization [17, 41] and self-training [27, 40, 49, 50, 52]. In this paper, we mainly focus

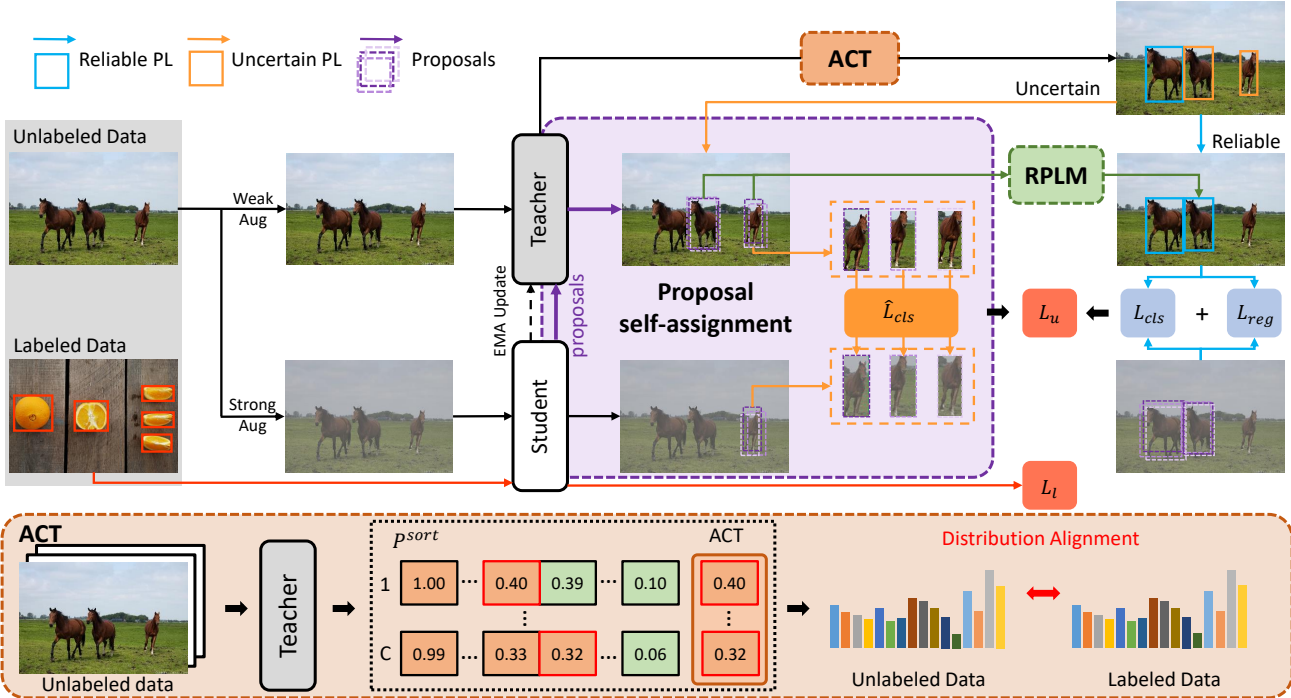


Figure 2. An overview of LabelMatch framework. Labeled data: only applied to the student with a supervised loss. Unlabeled data: annotated by the teacher to get pseudo labels (PL) according to the *adaptive label-distribution-aware confidence thresholds* (ACT), which are then split into reliable ones and uncertain ones for separated optimization. Reliable pseudo labels directly follow the IoU based assignment strategy, acting as hard labels to train the student model. As for uncertain labels, the *proposal self-assignment* method guides the student training with the supervision provided by the corresponding proposal prediction in the teacher. Besides, a *reliable pseudo label mining* (RPLM) strategy is utilized to convert the high-quality uncertain pseudo labels into reliable ones as the training goes on.

on the latter one. STAC [40] first generates pseudo labels by the pre-trained model and then feeds them back into the network with strong augmentation for model fine-tuning. To simplify this offline pseudo labeling, mean teacher based methods [27, 49, 52] perform a weak data transformation for online pseudo labeling and a strong data transformation for model training. However, there lies both foreground-background imbalance and foreground classes imbalance in SSOD, which makes it a more challenging task than semi-supervised classification. Unbiased-Teacher [27] and Soft-Teacher [49] use the focal-loss and the soft-weight to alleviate these problems, respectively. Despite the great progress, the label mismatch problem during pseudo labeling still exists in the previous works. In contrast, we propose a Label-Match framework to solve this problem from perspectives of distribution-level and instance-level.

Label Assignment. It is necessary to assign the target of classification and localization for each proposal or anchor in object detection, known as label assignment [11], which can be categorized into fixed and dynamic variants [11]. IoU based and center based label assignment are two common fixed assigning strategies, the first of which shows effectiveness in both RCNN-series [8, 12, 13, 29, 34] and one-stage detectors [24, 26, 32, 33], while the second one is popular

in many anchor-free object detection [20, 31, 43]. Recently, many adaptive mechanisms have been proposed to promote the label assignment, such as ATSS [51], PAA [19], AutoAssign [53], OTA [11], etc. However, all of these methods are only applied in supervised object detection, leaving a blank in SSOD due to the complex situation. To cope with the instance-level label mismatch problem, a novel label assignment is proposed to facilitate self-training in this paper.

3. Methodology

In SSOD, a set of labeled images $D_l = \{x_i^l, y_i^l\}_{i=1}^{N_l}$ and a set of unlabeled images $D_u = \{x_i^u\}_{i=1}^{N_u}$ are provided, where N_l and N_u represent the number of labeled and unlabeled data, respectively. The annotation y_i^l contains both categories and bounding boxes information.

3.1. Overview

The pipeline of LabelMatch framework is illustrated in Fig. 2, which is derived from a basic mean teacher framework. The main idea of the mean teacher framework is to drive the teacher and student to evolve in a mutual learning mechanism. However, previous mean teacher based works inevitably suffer from the label mismatch problems, which we divide into two granularities, including distribution-level

and instance-level. To solve these problems, we modify the mean teacher framework and develop a LabelMatch framework, consisting of a re-distribution mean teacher to solve the distribution-level label mismatch problem and a proposal self-assignment method to deal with the instance-level label mismatch problem. Also, it is beneficial to explore more high-quality pseudo labels. Thus, we further equip the proposed LabelMatch with a reliable pseudo label mining strategy to improve performance.

3.2. Preliminary: Mean Teacher Framework

Our approach follows the regimen of mean teacher, which contains a teacher model for pseudo label generation and a student model to improve the teacher model by updating knowledge. Both labeled and unlabeled data jointly constitute the batch of data. In each iteration, the teacher model first generates pseudo labels on the weakly-augmented unlabeled data, which are served as supervision signals for the corresponding strongly-augmented version. Subsequently, the student model is trained on the labeled data and the strongly-augmented unlabeled data with pseudo labels. In this way, the final training objective consists of a supervised loss and an unsupervised loss:

$$\mathcal{L}_l = \sum_i \mathcal{L}_{cls}(x_i^l, y_i^l) + \mathcal{L}_{reg}(x_i^l, y_i^l), \quad (1)$$

$$\mathcal{L}_u = \sum_i \mathcal{L}_{cls}(x_i^u, y_i^u) + \mathcal{L}_{reg}(x_i^u, y_i^u), \quad (2)$$

where \mathcal{L}_{cls} is the classification loss, \mathcal{L}_{reg} is the box regression loss, and y_i^u is the pseudo label annotated by the teacher model. The overall loss is defined as:

$$\mathcal{L}_{total} = \mathcal{L}_l + \lambda \mathcal{L}_u, \quad (3)$$

where λ is a weight to balance the unsupervised loss, which is set 2.0 by default in this paper. During self-training, the teacher gradually updates its weights from the student via an exponential moving average (EMA) strategy.

3.3. LabelMatch

We claim that the main obstacle to hinder the performance of mean teacher framework lies in the label mismatch problem. The proposed LabelMatch framework adopts the same mean teacher scheme, but develops a re-distribution mean teacher, which utilizes adaptive label distribution-aware confidence thresholds (ACT) to achieve unbiased pseudo labels. Moreover, a proposal self-assignment method and a reliable pseudo label mining strategy are introduced to rectify self-training.

Re-distribution Mean Teacher. In semi-supervised learning, the labeled data and the unlabeled data are from the same data distribution. Intuitively, we can obtain adaptive thresholds by minimizing the discrepancy of class distributions between the labeled and the unlabeled data, which can

be formulated as follows:

$$\begin{aligned} \underset{t_1, \dots, t_C}{\operatorname{argmin}} \quad & D_{KL}(\underbrace{[r_1^l, \dots, r_C^l]}_{f-f}, \underbrace{r_f^l}_{f-b}, \underbrace{[r_1^u, \dots, r_C^u]}_{f-f}, \underbrace{r_f^u}_{f-b}) \\ \text{s.t.} \quad & r_c^l = \frac{n_c^l}{\sum_{i=1}^C n_i^l}, \\ & r_f^l = \frac{\sum_{i=1}^C n_i^l}{N_l}, \\ & r_c^u = \frac{\sum_{j=1}^{N_u} \sum (P_i^j > t_c)}{\sum_{i=1}^C (\sum_{j=1}^{N_u} \sum (P_i^j > t_i))}, \\ & r_f^u = \frac{\sum_{i=1}^C (\sum_{j=1}^{N_u} \sum (P_i^j > t_i))}{N_u}, \end{aligned} \quad (4)$$

where D_{KL} represents the Kullback-Leibler divergence between two distributions, n_i^l denotes the box number of the i -th class in the labeled data, C is the entire foreground class number, P_i^j is a list of prediction scores of the i -th class in the j -th unlabeled image, $f-f$ means the foreground-foreground class distribution, and $f-b$ means the foreground-background ratio. Note that all the predictions in the unlabeled data are post-processed by NMS. t_c is the optimized variant, aka the confidence threshold to filter pseudo boxes for the c -th category, determined as:

$$t_c = P_c^{sort} \left[n_c^l \cdot \frac{N_u}{N_l} \right], \quad (5)$$

where P_c^{sort} is a list of prediction scores of the c -th class, which are sorted by descending. For efficient implementation, only a subset of unlabeled data are selected to estimate the distribution for thresholds determination. While the model is consecutively optimized during training, the previous thresholds are however imprecise for pseudo labeling, failing to be consistent with the truth class distribution. We thus simply update these thresholds every K iterations to dynamically adjust to the current teacher model. Such that, the thresholds are category-specific and adaptively up-to-date, termed as adaptive label-distribution-aware confidence thresholds (ACT), which we identify as the critical step to solve the distribution-level mismatch problem.

Proposal Self-Assignment. It is worth noting that the quality of pseudo labels cannot be guaranteed, especially at the early beginning of self-training. Inspired by the noise label learning [5, 14], we divide pseudo labels into reliable ones and uncertain ones according to the confidence score. Denoting $\alpha\%$ as the pre-defined proportion of reliable pseudo labels, the confidence thresholds t_c^r to filter reliable pseudo labels for the c -th category can be written as:

$$t_c^r = P_c^{sort} \left[\alpha\% \cdot n_c^l \cdot \frac{N_u}{N_l} \right], \quad (6)$$

Pseudo labels with confidence higher than t_c^r are regarded as hard labels for student model optimization in a supervised manner. In contrast, the remaining uncertain ones are treated as soft labels for soft learning.

Obviously, the uncertain pseudo labels potentially lead to low localization accuracy. To avoid poor box regression, Unbiased-Teacher removes the box regression loss for the unlabeled data, yet resulting in ambiguity in label assignment as shown in Fig. 1. For example, supposed that the proposals with an IoU overlap higher than 0.5 are optimized to the same uncertain pseudo labels, they will tend to be the same in classification score but different in localization after being refined by ROIHead. These refined proposals behave indistinguishably for the NMS post-processing, which confuses NMS to suppress redundant boxes randomly. We refer to this situation as an instance-level label mismatch problem, lacking attention in the previous SSOD works. To this end, we present a novel proposal self-assignment method for proposal re-calibration. Specifically, we utilize the proposals matched to the uncertain pseudo labels generated by the student to extract the corresponding features in the teacher, and then feed these features into the ROIHead of the teacher to achieve the refined boxes. Different from the IoU based label assignment, each proposal in the student uses the corresponding soft labels refined by the ROIHead of the teacher model for self-training, and finally, varying from each other in classification score to avoid NMS confusion. In this way, we optimize the student model with the uncertain pseudo labels via a soft classification loss:

$$\hat{\mathcal{L}}_{cls} = \sum_{i=1}^{n_p} \sum_{c=1}^C -p_{i,c}^t \log p_{i,c}^s, \quad (7)$$

where n_p is the number of the corresponding proposals matched to the uncertain pseudo labels, C is the class number, $p_{i,c}^s$ is the probability of the c -th class in the i -th proposal from the student model, and $p_{i,c}^t$ denotes the corresponding soft label from the teacher model matched to $p_{i,c}^s$.

Incorporating the re-distribution mean teacher and the proposal self-assignment for pseudo labeling, the unsupervised loss in Eq. (2) can be reformulated as:

$$\mathcal{L}_u = \sum_i \mathcal{L}_{cls}(x_i^u, y_i^{ur}) + \mathcal{L}_{reg}(x_i^u, y_i^{ur}) + \hat{\mathcal{L}}_{cls}(x_i^u, y_i^{uu}), \quad (8)$$

where y_i^{ur} and y_i^{uu} denote the reliable pseudo labels and the uncertain pseudo labels, respectively.

Reliable Pseudo Label Mining. To benefit from the continuously evolved teacher model and encourage the cycle positive feedback during self-training, we present a reliable pseudo label mining strategy to convert the high-quality uncertain pseudo labels into reliable ones. First, it is known that a set of adjacent boxes will be suppressed into one box after NMS. In this paper, we claim that these adjacent boxes before NMS can be exploited to evaluate the quality of the corresponding pseudo label after NMS. In this way, we present two evaluation metrics in this paper, *i.e.*, *mean score* and *mean IoU*, which are the average classification scores and IoU of this set of adjacent boxes matched to the corresponding bounding box after NMS. Note that We use the predictions from the teacher to compute *mean score* and

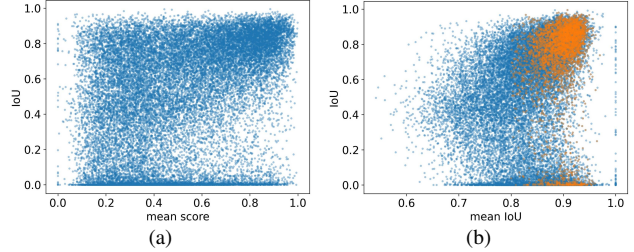


Figure 3. 5k images are selected to estimate the quality of pseudo labels. (a) the correlation between the IoU with ground truth and *mean score*. (b) the correlation between the IoU with ground truth and *mean IoU*. The orange points represent the predictions with the *mean score* larger than 0.8 and *mean IoU* larger than 0.8

mean IoU, and the IoU scores are determined by the IoU between the suppressed boxes and the selected box in NMS. We claim that the pseudo labels with higher quality usually correspond to higher mean scores and higher mean IoU. The empirical study in Fig. 3 gives a demonstration of our hypothesis. In this way, the uncertain pseudo labels with mean scores larger than T_{score} and mean IoU larger than T_{iou} will be transferred to reliable pseudo labels.

4. Experiments

4.1. Experimental Setup

Datasets. We evaluate our method on the MS-COCO [25] and PASCAL-VOC [10] datasets. There are three settings following the existing works [27, 40]: (1) COCO-standard: 1%, 5%, 10% images of train2017 set are sampled as the labeled training data and the remaining ones as the unlabeled data. (2) COCO-additional: we use the entire train2017 set as the labeled data and the additional COCO2017-unlabeled set as the unlabeled data. (3) VOC: we use VOC07 trainval set as the labeled data and the VOC12 trainval set as the unlabeled data. The validation sets in the COCO setting and VOC setting are COCO val2017 and VOC07 test set, respectively.

Network. For a fair comparison, we use Faster-RCNN [34] with FPN [23] and ResNet-50 backbone [15] as the detector. Our framework can be easily extended to other detectors.

Implementation Details. We implement our method with MMDetection [4]. For data augmentation, we apply random horizontal flipping and multi-scale for weak augmentation. Based on this augmentation, we then add random color jittering, grayscale, gaussian blurring and cutout patches for strong augmentation, which is similar to [27]. The T_{score} and T_{iou} in RPLM are set to 0.8 by default. More training and implementation details are introduced in the Appendix.

4.2. Results

COCO-standard. We evaluate the proposed method on COCO-standard (Tab. 1). Our method consistently outper-

	Threshold	1%	5%	10%
Supervised [27]	-	9.05 ± 0.16	18.47 ± 0.22	23.86 ± 0.81
STAC [40]	0.9	13.97 ± 0.35 (+4.92)	24.38 ± 0.12 (+5.91)	28.64 ± 0.21 (+4.78)
ISMT [50]	0.9	18.88 ± 0.74 (+9.83)	26.37 ± 0.24 (+7.90)	30.53 ± 0.52 (+6.67)
Instant Teaching [52]	0.9	18.05 ± 0.15 (+9.00)	26.75 ± 0.05 (+8.28)	30.40 ± 0.05 (+6.54)
Unbiased Teacher [27]	0.7	20.75 ± 0.12 (+11.70)	28.27 ± 0.11 (+9.80)	31.50 ± 0.10 (+7.64)
Soft Teacher [49]	0.9	20.46 ± 0.39 (+11.41)	30.74 ± 0.08 (+12.27)	34.04 ± 0.14 (+10.18)
LabelMatch (Ours)	ACT	25.81 ± 0.28 (+16.76)	32.70 ± 0.18 (+14.23)	35.49 ± 0.17 (+11.63)

Table 1. Experimental results on COCO-standard ($AP_{50:95}$). All the results are the average of all 5 folds.

	Iterations	$AP_{50:95}$
STAC [40]	540k	39.5 $\xrightarrow{-0.3}$ 39.2
Unbiased Teacher [27]	270k	40.2 $\xrightarrow{+1.1}$ 41.3
Soft Teacher [49]	370k	40.9 $\xrightarrow{+3.6}$ 44.5
LabelMatch (Ours)	540k	40.3 $\xrightarrow{+5.0}$ 45.3

Table 2. Experimental results on COCO-additional.

	AP_{50}	$AP_{50:95}$
Supervised [27]	72.63	42.13
STAC [40]	77.45 (+4.82)	44.64 (+2.51)
ISMT [50]	77.23 (+4.60)	46.23 (+4.10)
Instant Teaching [52]	79.20 (+6.57)	50.00 (+7.87)
Unbiased Teacher [27]	77.37 (+4.74)	48.69 (+6.56)
LabelMatch (Ours)	85.48 (+12.85)	55.11 (+12.98)

Table 3. Experimental results on VOC.

forms the previous state-of-the-arts with different percentages of labeled data. It is worth mentioning that the proposed method achieves 25.81 mAP on 1% labeled data, which is even higher than the supervised baseline trained on 10% labeled data.

COCO-additional. We verify whether the model trained on 100% COCO data can be further improved by using additional unlabeled COCO data. As shown in Tab. 2, our method boosts the supervised baseline with +5.0 mAP, while the existing SOTA improvement is +3.6 mAP.

VOC. We evaluate the proposed method on PASCAL-VOC. As is shown in Tab. 3, our method achieves 85.48 mAP on AP_{50} , which outperforms previous state-of-the-arts by +6.28 absolute mAP improvement.

4.3. Ablation Studies

In ablation studies, we conduct experiments in the setting of 1% COCO-standard (one of 5 folds), without RPLM strategy if not specified. The training iterations are reduced to 40k for all of the experiments. More implementation details and ablation studies can be found in the Appendix.

The quality of pseudo labels. The quality of pseudo la-

bels can be reflected in three aspects: 1) the accuracy of pseudo labels; 2) foreground-background distribution; 3) class distribution (foreground-foreground distribution). We compare the proposed method against the single confidence threshold based mean teacher framework. LabelMatch shows superior advancement, which we attribute to the following aspects:

- More accurate pseudo labels. As is shown in Fig. 4a, the accuracy of pseudo labels decreases when using threshold=0.7 and threshold=0.8. In contrast, the accuracy increases in LabelMatch. Although the accuracy achieves the best in threshold=0.9, the number (recall) of foregrounds is much less than the ground truth.
- Unbiased foregrounds-background distribution. As is shown in Fig. 4b, the number of pseudo labels in our method is nearly the same as the ground truth, while the number in the single confidence threshold based method is much lower, especially when threshold=0.9.
- Consistent class distribution. Fig. 4c demonstrates that LabelMatch guarantees both the foreground-background distribution and the class distribution to be nearly consistent with the ground truth. The situation is totally different when using the single confidence threshold, which brings large gaps with ground truth in many categories.

To show the quality of pseudo labels more intuitively, we give the quantitative and qualitative demonstrations in Tab. 4 and Fig. 5, respectively. From Tab. 4, we observe obvious gains in AP on top20 tail and head categories by leveraging LabelMatch compared with the single and fixed threshold. As shown in Fig. 5, there are many false positives with threshold=0.7 and many false negatives with threshold=0.9, which are removed by leveraging LabelMatch. These experimental results indicate the effectiveness of LabelMatch to re-distribute the pseudo label distribution, preventing self-training from collapsing to dominant classes. More qualitative results are shown in the Appendix.

The necessity of ACT adaptation. We randomly select a subset of unlabeled data to determine ACT for every K iterations in our implementation. As mentioned in Sec. 3.3, the proposed ACT are updated to the evolved teacher during the training phase, avoiding a false bias caused by the

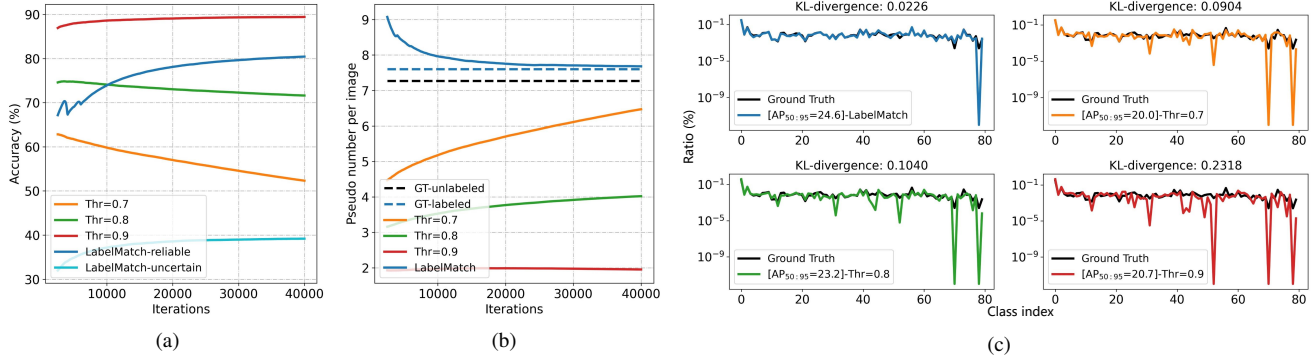


Figure 4. Ablation study about the quality of pseudo labels. (a) the accuracy of pseudo labels in the training phase. (Note: the pseudo labels with IoU overlapping the ground truth greater than 0.5 are regarded as true positives) (b) The average number of pseudo labels per image in the training phase. (c) KL divergence of the class distribution between the pseudo labels and the ground truth in the training phase.

	toaster	hair drier	scissors	microwave	toothbrush	parking meter	snowboard	bear	stop sign	fire hydrant	mous	frisbee	refrigerator	hot dog	oven	airplane	baseball bat	baseball glove	keyboard	bed	mean
thr=0.7	0.0	0.0	7.4	37.9	4.8	31.9	1.8	47.6	17.6	30.2	33.8	36.7	14.8	2.0	11.9	15.9	6.9	16.4	24.4	21.8	18.2
thr=0.8	0.0	0.0	7.7	37.2	4.7	32.0	3.0	51.5	32.4	46.2	44.8	44.5	33.7	4.6	18.5	34.8	11.1	24.0	21.8	27.1	24.0
thr=0.9	0.0	0.0	5.3	37.9	2.2	17.6	2.8	48.4	53.4	46.8	44.8	37.8	31.7	1.2	16.8	33.1	12.4	23.7	26.8	26.6	23.5
LabelMatch	4.6	0.0	16.5	37.5	4.6	34.9	9.5	49.9	50.4	49.2	47.4	45.0	34.4	8.5	17.7	39.3	13.2	21.0	29.3	32.3	27.3

	person	car	chair	book	bottle	cup	dining table	traffic light	bowl	handbag	bird	boat	truck	umbrella	bench	cow	banana	carrot	motorcycle	backpack	mean
thr=0.7	39.1	30.5	9.2	2.1	20.3	24.0	12.5	19.4	22.2	4.3	16.1	14.4	3.5	20.0	7.7	31.8	8.9	3.1	27.4	3.8	16.0
thr=0.8	39.0	31.8	10.3	2.1	25.8	26.3	14.1	20.2	25.6	4.1	18.1	13.6	8.6	20.0	11.4	32.1	8.3	3.9	28.0	4.8	17.4
thr=0.9	32.0	26.6	7.1	1.0	13.6	19.9	14.9	17.8	23.0	1.1	14.9	8.1	12.7	15.4	10.3	24.2	6.6	1.8	24.5	2.5	13.9
LabelMatch	38.8	31.9	11.7	1.9	25.6	26.1	14.6	20.8	29.2	4.8	17.9	13.2	14.8	21.2	13.6	39.0	9.1	4.9	27.4	6.1	18.6

Table 4. Quantitative results on top20 tail categories (upper) and top20 head categories (lower) ($AP_{50:95}$).

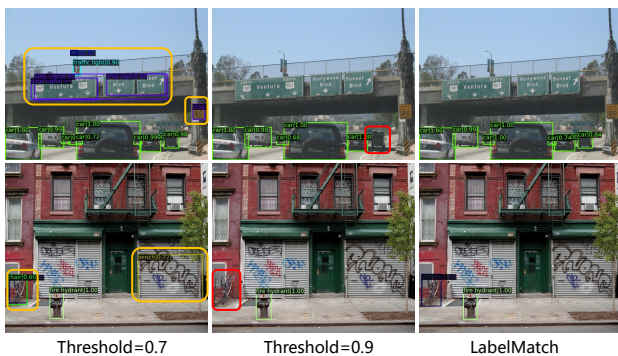


Figure 5. Qualitative comparisons between the single confidence threshold and our proposed method. Red rectangles highlight the false negatives, and yellow rectangles highlight the false positives. The score threshold for visualization is 0.6.

outdated predictions. Fig. 6a proves the necessity of ACT adaptation, where the performance is pretty bad without adaptation ($K = +\infty$). In our experiments, we simply use $K = 1000$ and a subset of unlabeled data (10,000) to refresh the thresholds.

Effect of reliable ratio α . In the training phase, we split

the candidate pseudo labels into reliable pseudo labels and uncertain pseudo labels by $\alpha\%$ according to the confidence scores. Here, we analyze the influence of different ratios. As shown in Fig. 6b, if we directly set all the candidate pseudo labels as uncertain pseudo labels ($\alpha = 0$), the performance is worse than splitting some pseudo labels to reliable, which is mainly caused by the lack of box regression optimization for the unlabeled data. However, setting too many pseudo labels as reliable pseudo labels is also harmful due to the noisy boxes. We use $\alpha = 20$ for all experiments.

Proposal self-assignment. We compare different label assignment strategies for uncertain pseudo labels: 1) Ignore assignment; 2) IoU based label assignment; 3) Proposal self-assignment. Here the ignore assignment means that the uncertain pseudo labels are directly set as ignore labels. Fig. 6c shows the superiority of the proposal self-assignment strategy over other label assignments. For the ignore assignment, there exists an imbalance between foreground and background, and background dominates the object detection training, which makes the ignored objects tend to be regarded as background after training, imped-

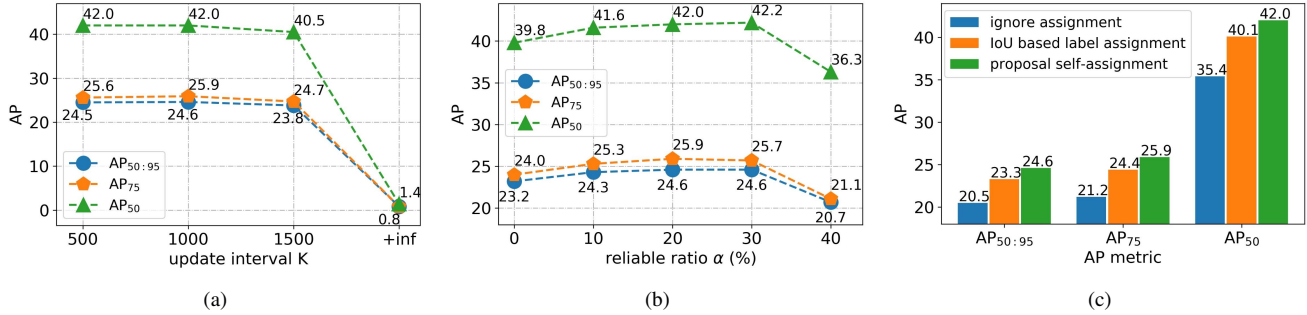


Figure 6. Ablation study about: (a) the updating interval of ACT. (b) the effect of reliable ratio α . (c) different label assignment strategies.

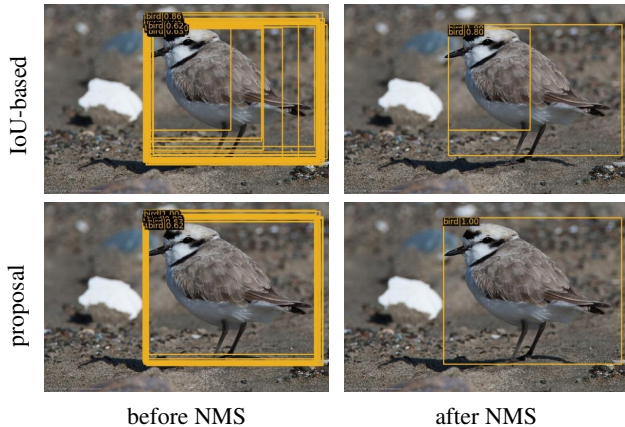


Figure 7. Visualization of the predictions before and after NMS post-processing. IoU based: the model is trained based on IoU based label assignment; proposal: the model is trained based on proposal self-assignment

ing performance improvement. For IoU based label assignment, it will produce many ambiguous boxes after training due to the instance mismatch as illustrated in Fig. 7.

Effect of RPLM. To verify the effectiveness of RPLM, we estimate it in the setting of COCO-standard. As depicted in Fig. 8a, RPLM slightly boosts the performance, but re-mitting the sensitivity of the reliable ratio α and showing an obvious improvement even α dropping to zero (Fig. 8b). This implies that, with the favor of RPLM, we can easily adjust α to make stable performance improvement without expert techniques.

5. Conclusion and Future Work

In this paper, we first diagnose the existing SSOD frameworks experimentally and figure out the common limitation, namely label mismatch problem, including two different but complementary granularities, *i.e.*, distribution-level and instance-level. To solve the above problems, we propose the LabelMatch framework. For distribution-level mismatch, LabelMatch develops a re-distribution mean teacher to derive adaptive label-distribution-aware confidence thresholds

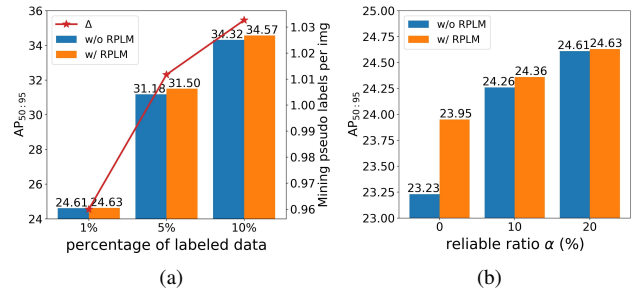


Figure 8. Ablation study about RPLM. (a) RPLM in different percentages of labeled data. (Δ means the number of mining pseudo labels per image.) (b) RPLM in different reliable ratio α .

by narrowing the class distribution discrepancy between labeled and unlabeled data, and then generating unbiased pseudo labels. For instance-level mismatch, LabelMatch adopts a proposal self-assignment method, injecting the proposals generated by the student into the teacher model to supervise proposals correction. Furthermore, a reliable pseudo label mining strategy is introduced to convert high-quality uncertain pseudo labels to reliable ones, facilitating the cycle of positive feedback during self-training. Extensive experimental results verify the efficacy of the proposed LabelMatch, which establishes a new state-of-the-art on both PASCAL-VOC and MS-COCO datasets.

Limitations. While we have shown the superiority of LabelMatch, there is still a non-negligible problem that the labeled and unlabeled data are assumed to follow the same distribution. Hence, LabelMatch relies on the class distribution prior, which is inaccessible in some scenarios, *e.g.*, unsupervised domain adaptive object detection. It is beneficial to study the label mismatch problems between two different distributions and, intuitively, advance the LabelMatch to more complex situations without class distribution prior, which is an interesting future work.

Acknowledgements

This work was sponsored in part by National Natural Science Foundation of China (U20B2066, 62106220), and Hikvision Open Fund.

References

- [1] Eric Arazo, Diego Ortego, Paul Albert, Noel E. O'Connor, and Kevin McGuinness. Pseudo-labeling and confirmation bias in deep semi-supervised learning. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8, 2020. **2**
- [2] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A. Raffel. Mixmatch: A holistic approach to semi-supervised learning. In *Advances in Neural Information Processing Systems*, volume 32, pages 5049–5059, 2019. **2**
- [3] Chaoqi Chen, Zebiao Zheng, Xinghao Ding, Yue Huang, and Qi Dou. Harmonizing transferability and discriminability for adapting object detectors. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8869–8878, 2020. **13**
- [4] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. Mmdetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019. **5**
- [5] Weijie Chen, LuoJun Lin, Shicai Yang, Di Xie, Shiliang Pu, Yueting Zhuang, and Wenqi Ren. Self-supervised noisy label learning for source-free unsupervised domain adaptation. *CoRR*, abs/2102.11614, 2021. **4**
- [6] Weijie Chen, Di Xie, Yuan Zhang, and Shiliang Pu. All you need is a few shifts: Designing efficient convolutional neural networks for image classification. In *CVPR*, 2019. **1**
- [7] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3213–3223, 2016. **13**
- [8] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-fcn: Object detection via region-based fully convolutional networks. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, volume 29, pages 379–387, 2016. **3**
- [9] Jinhong Deng, Wen Li, Yuhua Chen, and Lixin Duan. Unbiased mean teacher for cross-domain object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4091–4101, 2021. **13**
- [10] Mark Everingham, Luc Gool, Christopher K. Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88(2):303–338, 2010. **2, 5**
- [11] Zheng Ge, Songtao Liu, Zeming Li, Osamu Yoshie, and Jian Sun. Ota: Optimal transport assignment for object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 303–312, 2021. **3**
- [12] Ross Girshick. Fast r-cnn. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 1440–1448, 2015. **3**
- [13] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR '14 Proceedings of the 2014 IEEE Conference on Computer Vision and Pattern Recognition*, pages 580–587, 2014. **3**
- [14] Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor W. Tsang, and Masashi Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels. In *Advances in Neural Information Processing Systems*, volume 31, pages 8527–8537, 2018. **4**
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. **5, 13**
- [16] Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry P. Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization. In *34th Conference on Uncertainty in Artificial Intelligence 2018, UAI 2018*, pages 876–885, 2018. **2**
- [17] Jisoo Jeong, Seungeui Lee, Jeesoo Kim, and Nojun Kwak. Consistency-based semi-supervised learning for object detection. In *Advances in Neural Information Processing Systems*, volume 32, pages 10758–10767, 2019. **1, 2**
- [18] Matthew Johnson-Roberson, Charles Barto, Rounak Mehta, Sharath Nittur Sridhar, Karl Rosaen, and Ram Vasudevan. Driving in the matrix: Can virtual worlds replace human-generated annotations for real world tasks? In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 746–753, 2017. **13**
- [19] Kang Kim and Hee Seok Lee. Probabilistic anchor assignment with iou prediction for object detection. In *European Conference on Computer Vision*, pages 355–371, 2020. **3**
- [20] Tao Kong, Fuchun Sun, Huaping Liu, Yuning Jiang, Lei Li, and Jianbo Shi. Foveabox: Beyond anchor-based object detection. *IEEE Transactions on Image Processing*, 29:7389–7398, 2020. **3**
- [21] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. *International Journal of Computer Vision*, 128(3):642–656, 2020. **1**
- [22] Xianfeng Li, Weijie Chen, Di Xie, Shicai Yang, Peng Yuan, Shiliang Pu, and Yueting Zhuang. A free lunch for unsupervised domain adaptive object detection without source data. In *AAAI*, 2021. **13**
- [23] Tsung-Yi Lin, Piotr Dollar, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 936–944, 2017. **5, 13**
- [24] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollar. Focal loss for dense object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(2):318–327, 2020. **3**
- [25] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755, 2014. **2, 5**

- [26] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott E. Reed, Cheng-Yang Fu, and Alexander C. Berg. Ssd: Single shot multibox detector. In *14th European Conference on Computer Vision, ECCV 2016*, pages 21–37, 2016. [3](#)
- [27] Yen-Cheng Liu, Chih-Yao Ma, Zijian He, Chia-Wen Kuo, Kan Chen, Peizhao Zhang, Bichen Wu, Zsolt Kira, and Peter Vajda. Unbiased teacher for semi-supervised object detection. In *ICLR 2021: The Ninth International Conference on Learning Representations*, 2021. [1](#), [2](#), [3](#), [5](#), [6](#), [14](#), [15](#)
- [28] Takeru Miyato, Shin-Ichi Maeda, Masanori Koyama, and Shin Ishii. Virtual adversarial training: A regularization method for supervised and semi-supervised learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(8):1979–1993, 2019. [2](#)
- [29] Jiangmiao Pang, Kai Chen, Jianping Shi, Huajun Feng, Wanli Ouyang, and Dahua Lin. Libra r-cnn: Towards balanced learning for object detection. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 821–830, 2019. [3](#)
- [30] Rindra Ramamonjison, Amin Banitalebi-Dehkordi, Xinyu Kang, Xiaolong Bai, and Yong Zhang. Simrod: A simple adaptation method for robust object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3570–3579, 2021. [14](#)
- [31] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 779–788, 2016. [3](#)
- [32] Joseph Redmon and Ali Farhadi. Yolo9000: Better, faster, stronger. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6517–6525, 2017. [1](#), [3](#)
- [33] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018. [3](#)
- [34] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(6):1137–1149, 2017. [1](#), [2](#), [3](#), [5](#)
- [35] Kuniaki Saito, Yoshitaka Ushiku, Tatsuya Harada, and Kate Saenko. Strong-weak distribution alignment for adaptive object detection. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6956–6965, 2019. [13](#), [14](#)
- [36] Mehdi Sajjadi, Mehran Javanmardi, and Tolga Tasdizen. Regularization with stochastic transformations and perturbations for deep semi-supervised learning. In *NIPS’16 Proceedings of the 30th International Conference on Neural Information Processing Systems*, volume 29, pages 1171–1179, 2016. [2](#)
- [37] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Semantic foggy scene understanding with synthetic data. *International Journal of Computer Vision*, 126(9):973–992, 2018. [13](#)
- [38] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *Computer Science*, 2014. [13](#)
- [39] Kihyuk Sohn, David Berthelot, Chun-Liang Li, Zizhao Zhang, Nicholas Carlini, Ekin D. Cubuk, Alex Kurakin, Han Zhang, and Colin Raffel. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. In *Advances in Neural Information Processing Systems*, volume 33, pages 596–608, 2020. [2](#)
- [40] Kihyuk Sohn, Zizhao Zhang, Chun-Liang Li, Han Zhang, Chen-Yu Lee, and Tomas Pfister. A simple semi-supervised learning framework for object detection. *arXiv preprint arXiv:2005.04757*, 2020. [1](#), [2](#), [3](#), [5](#), [6](#), [11](#), [14](#)
- [41] Peng Tang, Chetan Ramaiah, Yan Wang, Ran Xu, and Caiming Xiong. Proposal learning for semi-supervised object detection. In *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 2290–2300, 2021. [2](#)
- [42] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *ICLR (Workshop)*, 2017. [2](#)
- [43] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9627–9636, 2019. [1](#), [3](#)
- [44] Vibashan Vs, Vikram Gupta, Poojan Oza, Vishwanath A. Sindagi, and Vishal M. Patel. Mega-cda: Memory guided attention for category-aware unsupervised domain adaptive object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4516–4526, 2021. [13](#), [14](#)
- [45] Qizhe Xie, Zihang Dai, Eduard H. Hovy, Thang Luong, and Quoc Le. Unsupervised data augmentation for consistency training. In *Advances in Neural Information Processing Systems*, volume 33, pages 6256–6268, 2020. [2](#)
- [46] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V. Le. Self-training with noisy student improves imagenet classification. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10687–10698, 2020. [2](#)
- [47] Chang-Dong Xu, Xing-Ran Zhao, Xin Jin, and Xiu-Shen Wei. Exploring categorical regularization for domain adaptive object detection. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11724–11733, 2020. [13](#)
- [48] Minghao Xu, Hang Wang, Bingbing Ni, Qi Tian, and Wenjun Zhang. Cross-domain detection via graph-induced prototype alignment. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12355–12364, 2020. [13](#), [14](#)
- [49] Mengde Xu, Zheng Zhang, Han Hu, Jianfeng Wang, Lijuan Wang, Fangyun Wei, Xiang Bai, and Zicheng Liu. End-to-end semi-supervised object detection with soft teacher. *arXiv preprint arXiv:2106.09018*, 2021. [1](#), [2](#), [3](#), [6](#), [14](#)
- [50] Qize Yang, Xihan Wei, Biao Wang, Xian-Sheng Hua, and Lei Zhang. Interactive self-training with mean teachers for semi-supervised object detection. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 5941–5950, 2021. [2](#), [6](#)
- [51] Shifeng Zhang, Cheng Chi, Yongqiang Yao, Zhen Lei, and Stan Z. Li. Bridging the gap between anchor-based and anchor-free detection via adaptive training sample selection.

In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9759–9768, 2020. 3

- [52] Qiang Zhou, Chaohui Yu, Zhibin Wang, Qi Qian, and Hao Li. Instant-teaching: An end-to-end semi-supervised object detection framework. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4081–4090, 2021. 1, 2, 3, 6
- [53] Benjin Zhu, Jianfeng Wang, Zhengkai Jiang, Fuhang Zong, Songtao Liu, Zeming Li, and Jian Sun. Autoassign: Differentiable label assignment for dense object detection. *arXiv preprint arXiv:2007.03496*, 2020. 3

A. Consistent Class Distribution Assumption

LabelMatch is based on the assumption that consistent class distribution exists between the labeled and unlabeled data since they are drawn from the same data distribution. To further verify this hypothesis, we present the comparisons between the labeled and unlabeled data in COCO-standard and VOC using the ground-truth labels. As shown in Fig. 9, the foreground-foreground class distribution and the foreground-background ratio of the unlabeled data are close to those of the labeled data in these SSOD settings.

B. More Results on COCO-standard

In this section, we present more experimental results on COCO-standard using the ablation study setting (see the fifth column in Tab. 13). Firstly, we carry out more analysis about ACT in Appendix B.1. Then, we study the effect of hyper-parameter in RPLM in Appendix B.2 and more analysis about proposal self-assignment in Appendix B.3. Finally, more qualitative results are exhibited in Appendix B.4.

B.1. Analysis of ACT

In this part, we present more analysis about the proposed ACT from flexibility and implementation.

Flexibility. To further demonstrate the flexibility of our method, we extend STAC [40] with the proposed ACT, denoted as STAC* for short. The original STAC first uses a pretrained model to generate pseudo labels and then uses a threshold of 0.9 to filter out low-quality pseudo labels, which are finally fed back into the network with strong data augmentation for model fine-tuning. Alternatively, STAC* replaces the fixed threshold with the proposed ACT for pseudo labeling and updates the thresholds every epoch. Since there is no mean teacher in STAC, the label assignment strategy of STAC* simply follows the *ignore assignment*, where uncertain pseudo labels are set as ignore labels. As shown in Fig. 10, there is an apparent performance gain after equipping STAC with ACT, demonstrating the universality of the proposed ACT.

Online vs. Offline. As discussed in the paper, ACT are updated to the evolved teacher during the training phase,

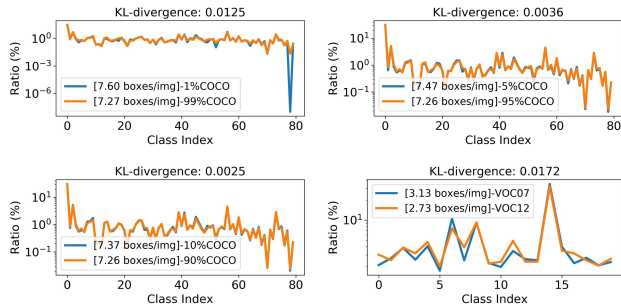


Figure 9. Comparisons on class distribution between the labeled and unlabeled data. The blue and orange lines denote the foreground-foreground class distribution in the labeled and unlabeled data, respectively. “boxes/img” in the legend represents the foreground-background ratio.

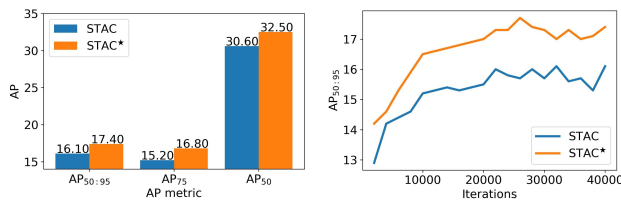


Figure 10. Performance comparisons between STAC and STAC* on COCO-standard with 1% labeled data.

	ACT	iterations	1%	5%	10%
LabelMatch	Offline	40K	24.6	31.6	34.6
LabelMatch	Online	40K	24.6	31.5	34.6

Table 5. Performance ($AP_{50:95}$) comparisons between the online and offline versions of ACT. We only run 1-fold using the ablation training setting.

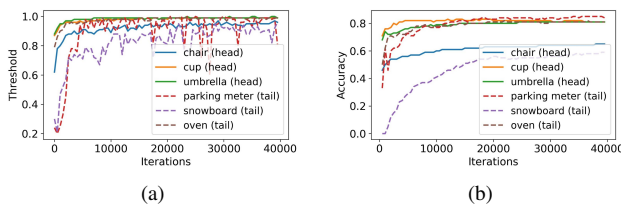


Figure 11. (a) Thresholds in the training phase. (b) The quality of reliable pseudo labels.

avoiding a negative bias caused by the outdated predictions. There are two patterns to update ACT, one of which is introduced in the paper, leveraging a subset of unlabeled data to update ACT every K iterations, termed as offline version. Here, we describe another pattern, named as online version, which maintains a scores queue, as shown in Fig. 12a. The teacher’s prediction is pushed into the scores queue for refreshing the ACT in each training iteration, which can be seen as a special case of the offline version with $K = 1$.

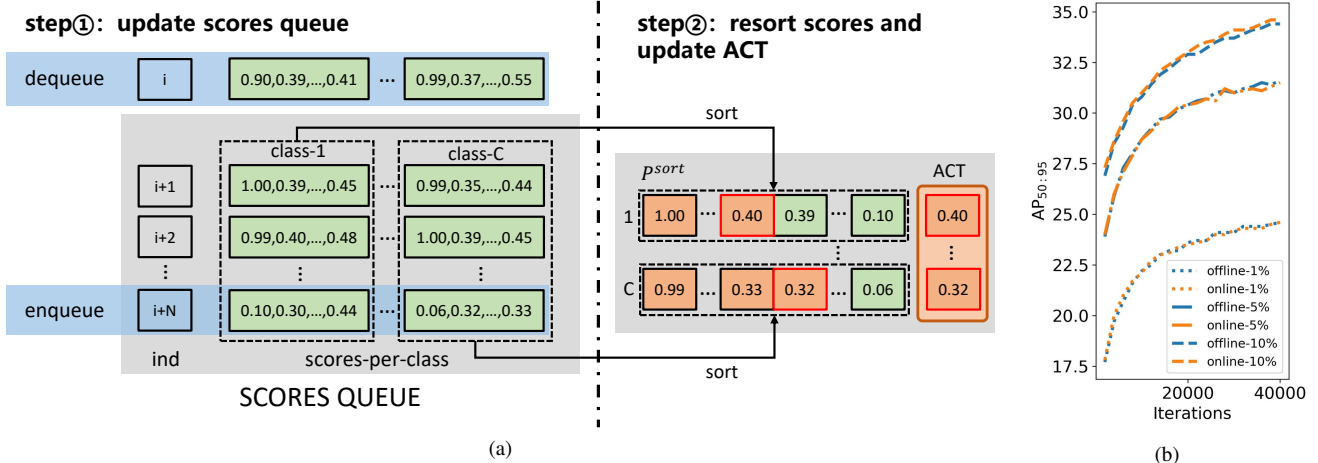


Figure 12. (a) Implementation of the online version of ACT. Each training iteration consists of two steps: (1) obtain predictions from the teacher model and update the scores queue; (2) re-sort scores and update the ACT in real-time. (b) Performance comparison between the online and offline version of ACT during the training phase on three COCO-standard settings with 1%, 5% and 10% labeled data.

Both versions of ACT can get satisfactory performance, as shown in Fig. 12b and Tab. 5. We use the offline version in all the experiments and will release the online version as well.

Thresholds evolve alone training. We select three head classes and three tail classes on offline version for analysis. As shown in Fig. 11, the thresholds (Eq.6 in the paper to filter reliable pseudo labels) increase in both head and tail classes during optimization. Specifically, the thresholds for tail classes are more fluctuated than those for head classes due to the scarce samples. Also, the quality of reliable pseudo labels get increased as training goes on.

B.2. Analysis of RPLM Hyper-Parameter

There are two hyper-parameters (T_{score}, T_{iou}) in the component of reliable pseudo label mining (RPLM). Here we use COCO-standard with 10% labeled data as the experimental setting. As shown in Tab. 6, the best performance appears when $(T_{score}, T_{iou}) = (0.8, 0.8)$. Therefore, we use $(T_{score}, T_{iou}) = (0.8, 0.8)$ by default in all experiments throughout the paper. It is also worth mentioning that our method is not sensitive to these hyper-parameters.

B.3. Analysis of proposal self-assignment

To further analyze the quality of the teacher’s RoI head predictions on the student’s proposals (proposals self-assignment vs. IoU-based label assignment), we use the ground truth for quantitative measurement. For each proposal, we calculate the cross-entropy between the corresponding prediction and the nearest ground truth (set as background if $\text{IoU} < 0.5$). As shown in Fig. 13, the predictions by proposal self-assignment show better quality than IoU-based one.

(T_{score}, T_{iou})	$AP_{50:95}$	(T_{score}, T_{iou})	$AP_{50:95}$	(T_{score}, T_{iou})	$AP_{50:95}$
(0.7, 0.7)	34.4	(0.8, 0.7)	34.4	(0.9, 0.7)	34.4
(0.7, 0.8)	34.5	(0.8, 0.8)	34.6	(0.9, 0.8)	34.5
(0.7, 0.9)	34.5	(0.8, 0.9)	34.6	(0.9, 0.9)	34.3

Table 6. Effect of hyper-parameters in RPLM.

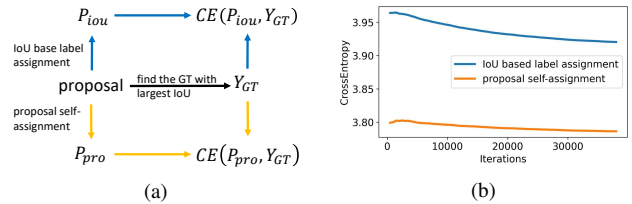


Figure 13. (a) The solution to measure the quality of predictions. (b) The quality comparison of the predictions between two label assignment methods in the training phase (lower is better).

B.4. Qualitative Results

We perform the qualitative comparisons between the proposed method and the mean teacher frameworks with a fixed and single confidence threshold (varying from 0.7 to 0.9). As shown in Fig. 14, there are many false positives with a low confidence threshold (yellow rectangles in the second column), while many false negatives appear when using a high confidence threshold (red rectangles in the fourth column). Although the manual search threshold (0.8) via trial-and-error can achieve satisfactory results, our method shows even better qualitative results.

C. Domain Adaptive Object Detection

LabelMatch is based on the consistent class distribution assumption between the labeled and unlabeled data. To ex-



Figure 14. Qualitative comparisons between the single confidence threshold and the proposed LabelMatch. Red rectangles highlight the false negatives, and yellow rectangles highlight the false positives. The score threshold for visualization is 0.6.

Data Split	Normal→Foggy	Small→Large	Across cameras	Synthetic→Real
labeled data	Cityscapes (train)	Cityscapes (train)	KITTI	Sim10K
unlabeled data	Cityscapes-foggy (train)	BDD100K (train)	Cityscapes (train)	Cityscapes (train)
test data	Cityscapes-foggy (val)	BDD100K (val)	Cityscapes (val)	Cityscapes (val)

Table 7. Four different domain shifts in DA-OD, which are constructed by five different datasets, including Cityscapes [7], Cityscapes-foggy [37], KITTI [37], Sim10k [18] and BDD100K [18].

Method	truck	car	rider	person	train	motor	bicycle	bus	mean
Source only	19.2	47.9	40.8	34.8	7.8	24.2	36.0	36.4	30.9
CVPR2020:GPA [48]	24.7	54.1	46.7	32.9	41.1	32.4	38.7	45.7	39.5
CVPR2020:HTCN [3]	31.6	47.9	47.5	33.2	40.9	32.3	37.1	47.4	39.8
CVPR2021:MeGA [44]	25.4	52.4	49.0	37.7	46.9	34.5	39.0	49.2	41.8
CVPR2021:UMT [9]	34.1	48.6	46.7	33.0	46.8	30.4	37.3	56.5	41.7
LabelMatch (Ours)	42.0	62.2	55.4	45.3	55.1	43.5	51.5	64.1	52.4

Table 8. Results of adaptation from normal to foggy weathers. “Source only” refers to the model trained by labeled source data.

explore the robustness of LabelMatch to the prior dependence on this assumption, we extend it to the scenario of domain adaptive object detection (DA-OD) [22,35,44,48] where the labeled source data and the unlabeled target are drawn from two different data distributions.

Dataset. As described in Tab. 7, following the existing DA-OD works, there are four common types of domain shifts in DA-OD. We evaluate our method on these settings and compare it with the state-of-the-arts.

Network Architecture. For a fair comparison with the existing DA-OD arts, we switch the backbone from ResNet-

Method	truck	car	rider	person	train	motor	bicycle	bus	mean
Source only	18.3	50.0	33.3	35.8	-	18.4	27.6	17.0	28.7
CVPR2019:SW-Faster [35]	15.2	45.7	29.5	30.2	-	17.1	21.2	18.4	25.3
CVPR2020:CR-DA [47]	19.5	46.3	31.3	31.4	-	17.3	23.8	18.9	26.9
LabelMatch (Ours)	39.4	54.6	37.4	42.9	-	25.7	29.8	41.7	38.8
LabelMatch [†] (Ours)	39.8	55.4	44.5	44.8	-	38.6	41.5	47.1	44.5

Table 9. Results of adaptation from small to large scale datasets. [†] is an ideal setting that uses the ground-truth labels of the unlabeled data for class distribution estimation.

50 [15] to VGG-16 [38] and remove the FPN [23] neck.

Implementation Details. The implementation is nearly the same as SSOD, and more training hyper-parameters can be found in Appendix E. Following previous works, we use AP_{50} as our evaluation metric.

Results. To examine the prior dependence on the consistent class distribution assumption, we evaluate LabelMatch in two class distribution estimation manners: 1) The first one is the same as described in the main body of the paper, which estimates the class distribution of the unlabeled target data by the annotations of the labeled source data; 2) The

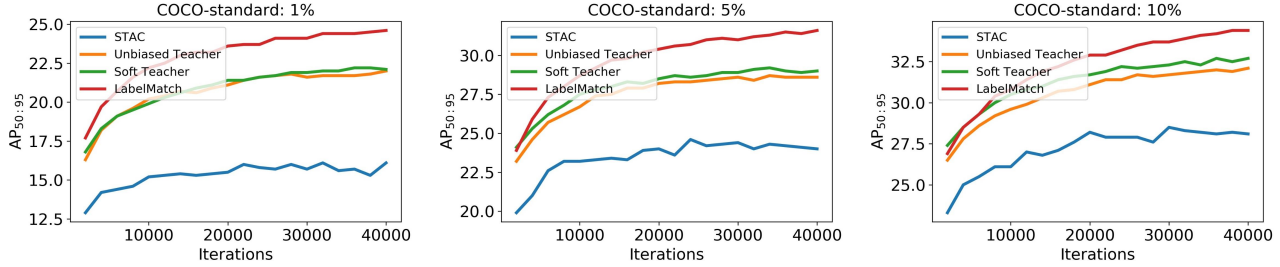


Figure 15. Performance ($AP_{50:95}$) comparisons among different state-of-the-art SSOD methods with exactly the same training settings.

Method	AP_{50}	network	Method	AP_{50}	network
Source only	42.2	FR+VGG	Source only	36.5	FR+VGG
CVPR2019:SW-Faster [35]	37.9	FR+VGG	CVPR2019:SW-Faster [35]	40.7	FR+VGG
CVPR2020:GPA [48]	47.9	FR+R50	CVPR2020:GPA [48]	47.6	FR+R50
CVPR2021:MeGA [44]	43.0	FR+VGG	CVPR2021:MeGA [44]	44.8	FR+VGG
ICCV2021:SimROD [30]	47.5	YOLOv5	ICCV2021:SimROD [30]	52.1	YOLOv5
LabelMatch (Ours)	51.0	FR+VGG	LabelMatch (Ours)	52.7	FR+VGG
LabelMatch [†] (Ours)	52.2	FR+VGG	LabelMatch [†] (Ours)	53.8	FR+VGG

Table 10. Results of adaptation across cameras. FR: Faster-RCNN. [†] is an ideal setting that uses the ground-truth labels of the unlabeled data for class distribution estimation.

Table 11. Results of adaptation from synthetic to real. VGG: VGG-16. [†] is an ideal setting that uses the ground-truth labels of the unlabeled data for class distribution estimation.

second one is an ideal setting, which determines the class distribution of the unlabeled target data by the ground-truth labels of the unlabeled data.

- **Normal→Foggy:** This scenario is different from the following DA-OD settings. In this scenario, the labeled source data and the unlabeled target data meet exactly the same class distribution since the target foggy data is rendered from the normal source data via a foggy translation model. As shown in Tab. 8, benefited from the given class distribution, we achieve a +21.5 mAP improvement over the “source only” baseline, exceeding previous state-of-the-arts by a large margin.
- **Small→Large:** Although there exists bias between the labeled and unlabeled data on foreground-foreground class distribution ($KL = 0.36$) and foreground-background ratio (18.5 boxes/img vs. 13.9 boxes/img), our method can still achieve 38.8 mAP, surpassing all the previous arts as far as we know. With access to the accurate class distribution (the ideal setting), our method can be further improved to 44.5 mAP.
- **Across cameras & Synthetic→Real:** In these settings, there is only one foreground class and exists foreground-background ratio bias (4.3 boxes/img vs. 9.6 boxes/img and 5.8 boxes/img vs. 9.6 boxes/img). Even using a biased class distribution, our method can still achieve satisfactory results. And our method can get further improve-

Method	Loss	Threshold	1%	5%	10%
STAC [40]	Cross-Entropy	0.9	16.1	24.0	28.1
Unbiased Teacher [27]	Focal-Loss	0.7	22.0	28.6	32.1
Soft Teacher* [49]	Cross-Entropy	0.9	22.1	29.0	32.7
LabelMatch (Ours)	Cross-Entropy	ACT	24.6	31.5	34.6

Table 12. Benchmark results on COCO-standard: our re-implementations with exactly the same training details and data augmentation strategies. * denotes the re-implementation without box-jitter trick. We only run 1-fold using the ablation training setting due to the limitation of computation resources.

ment equipped with the accurate class distribution (aka the ideal setting).

These DA-OD experiments demonstrate the robustness of the proposed LabelMatch framework, since the introduction of proposal self-assignment and RPLM weaken the prior dependence on the consistent class distribution assumption. From another perspective, these experiments also indicate that an accurate class distribution estimation can further promote the performance of DA-OD, emphasizing the importance of class distribution estimation. How to estimate an accurate class distribution when the labeled data and the unlabeled data are drawn from two different data distributions is an interesting future work.

D. MMDetection-based SSOD Codebase

Since different SSOD algorithms use different data augmentation strategies which have great impact on the performance, we build a unified MMDetection-based SSOD codebase for a fair comparison, named MMDet-SSOD for short, containing STAC [40], Unbiased-Teacher [27], Soft-Teacher [49] and LabelMatch.

We comprehensively run all algorithms in our MMDet-SSOD on COCO-standard dataset using the ablation training setting, and report the performance in Tab. 12 and Fig. 15. It is worth mentioning that the data augmentation, training iterations, batch size, and other training settings are all kept the same among these algorithms for a fair comparison. The entire source code will be released soon to support the development of SSOD in the community.

training setting	COCO-standard	COCO-additional	VOC	Ablation	DA-OD
batch size for labeled data	16	32	4	32	16
batch size for unlabeled data	16	32	4	32	16
learning rate	0.01	0.02	1.25e-3	0.02	0.016
learning rate step	-	(360K, 480K)	-	-	-
iterations	160K	540K	160K	40K	20K
unsupervised loss weight λ	2.0	2.0	2.0	2.0	2.0
EMA rate	0.996	0.996	0.996	0.996	0.9996
reliable ratio α	0.2	0.2	0.2	0.2	0.2
mean score thresh T_{score}	0.8	0.8	0.8	0.8	0.8
mean iou thresh T_{iou}	0.8	0.8	0.8	0.8	0.8
multi-scale (strong augmentation)	(0.2, 1.8)	(0.2, 1.8)	(0.2, 1.8)	(0.5, 1.5)	(0.5, 1.5)
test score thresh	0.001	0.001	0.001	0.001	0.001

Table 13. Training settings for different datasets and different tasks. ‘‘Ablation’’ means the training setting of the ablation studies in the main body of the paper, which is also used in all SSOD experiments in the Appendix.

Weak Augmentation			
Process	Prob	Parameters	Descriptions
Horizontal Flip	0.5	None	None
Multi-Scale	1.0	scale=(500, 800)	The short edge of image is random resized from 500 to 800.
Strong Augmentation			
Process	Prob	Parameters	Descriptions
Horizontal Flip	0.5	None	None
Multi-Scale	1.0	ratio=(0.2, 1.8)	The short edge of image is random resized from $0.5l_{short}$ to $1.5l_{short}$.
Color Jittering	0.8	(brightness, contrast, saturation, hue) = (0.4, 0.4, 0.4, 0.1)	Brightness factor is chosen uniformly form [0.6, 1.4], contrast factor is chosen uniformly from [0.6, 1.4], saturation factor is chosen uniformly from [0.6, 1.4], and hue value is chosen uniformly from [-0.1, 0.1].
Grayscale	0.2	None	None
GaussianBlur	0.5	(sigma_x, sigma_y)=(0.1, 2.0)	Gaussian filter with $\sigma_x = 0.1$ and $\sigma_y = 2.0$ is applied
CutoutPattern1	0.7	scale=(0.05, 0.2), ratio=(0.3, 3.3)	Randomly selects a rectangle region in an image and erases its pixels.
CutoutPattern2	0.7	scale=(0.02, 0.2), ratio=(0.1, 6.0)	Randomly selects a rectangle region in an image and erases its pixels.
CutoutPattern3	0.7	scale=(0.02, 0.2), ratio=(0.05, 8.0)	Randomly selects a rectangle region in an image and erases its pixels.

Table 14. Details of data augmentations. In our ablation study, we use multi-scale with ratio=(0.5, 1.5) in order to use large batch size.

E. Implementation and Training Details

Training. We utilize different training settings for different datasets in our implementation. We use the SGD optimizer with a momentum rate 0.9 and weight decay 0.0001 in all experiments. The different training settings are summarized in Tab. 13.

Data augmentation. Our data augmentation strategies are modified from Unbiased Teacher [27], and the details are shown in Tab. 14. The weak augmentation is applied to the unlabeled data for pseudo labeling, and the strong augmentation is applied to both labeled and unlabeled data for model training. In our implementation, no cutout augmentation is applied to the labeled data when using strong data

augmentation. In order to save computation resources, we use multi-scale with ratio=(0.5, 1.5) in the ablation studies.