

# Efficient Teacher: Semi-Supervised Object Detection for YOLOv5

Bowen Xu   Mingtao Chen   Wenlong Guan   Lulu Hu  
Alibaba Group

{bowen.xbw, ruiyang.cmt, wenlong.gwl, chudu.hll}@alibaba-inc.com

## Abstract

*Semi-Supervised Object Detection (SSOD) has been successful in improving the performance of both R-CNN series and anchor-free detectors. However, one-stage anchor-based detectors lack the structure to generate high-quality or flexible pseudo labels, leading to serious inconsistency problems in SSOD. In this paper, we propose the Efficient Teacher framework for scalable and effective one-stage anchor-based SSOD training, consisting of Dense Detector, Pseudo Label Assigner, and Epoch Adaptor. Dense Detector is a baseline model that extends RetinaNet with dense sampling techniques inspired by YOLOv5. The Efficient Teacher framework introduces a novel pseudo label assignment mechanism, named Pseudo Label Assigner, which makes more refined use of pseudo labels from Dense Detector. Epoch Adaptor is a method that enables a stable and efficient end-to-end semi-supervised training schedule for Dense Detector. The Pseudo Label Assigner prevents the occurrence of bias caused by a large number of low-quality pseudo labels that may interfere with the Dense Detector during the student-teacher mutual learning mechanism, and the Epoch Adaptor utilizes domain and distribution adaptation to allow Dense Detector to learn globally distributed consistent features, making the training independent of the proportion of labeled data. Our experiments show that the Efficient Teacher framework achieves state-of-the-art results on VOC, COCO-standard, and COCO-additional using fewer FLOPs than previous methods. To the best of our knowledge, this is the first attempt to apply Semi-Supervised Object Detection to YOLOv5.*

## 1. Introduction

Object detection [3, 25, 31, 40] has made significant advances in recent years, which follows a traditional supervised training approach and relies on costly manual annotation efforts. To mitigate this problem, many semi-supervised techniques [1, 35] are proposed to exploit large amounts of unlabeled data by automatically generating pseudo labels without introducing manual annotation. De-

spite great progress in SSOD [4, 5, 27, 45], there are three key issues that remain challenging:

Firstly, **few works on one-stage anchor-based SSOD have been reported**. Though anchor-free detectors [15, 22, 40] have been recently getting more attention in the community of object detection, one-stage anchor-based detectors [2, 19, 22, 30, 42], having the advantages of high recall, high numerical stability and fast training speed, are widely used in scenarios with extremely high recall demands. However, most SSOD methods are implemented on a two-stage anchor-based detector such as Faster R-CNN [31] and an one-stage anchor-free detector such as FCOS [40], which output relatively sparse bounding box predictions due to the multi-stage coarse-to-fine prediction mechanism or the anchor-free design of detection head. In contrast, the classic one-stage anchor-based detector generates more dense predictions due to its multiple-anchor mechanism, which leads to positive and negative samples imbalance during supervised training [46] and poor quality of pseudo labels during semi-supervised training.

Secondly, current mainstream SSOD approaches, following a student-teacher mutual learning manner [27, 45], is difficult for an one-stage anchor-based detector to train due to the serious **pseudo label inconsistency** problem, that is, throughout the training process, the quantity and quality of pseudo labels generated by the teacher model fluctuates greatly and the unqualified pseudo labels can mislead model updates. To alleviate this problem, two-stage methods [4] [27] refine pseudo labels several times more than one-stage methods and anchor-free methods [49] adopt feature maps as soft pseudo labels to avoid bias caused by non maximum suppression. The pseudo label inconsistency is exacerbated in an one-stage anchor-based detector because of its multiple-anchor mechanism mentioned above. The work [47] has reported that the SSOD experimental results of RetinaNet are not as good as those on Faster R-CNN and FCOS.

Thirdly, **how to train a SSOD model with both higher accuracy and better efficiency** becomes the key issue that restricts the application of SSOD in a wide range of scenarios. The previous SSOD methods [4, 23, 27, 45, 50] are

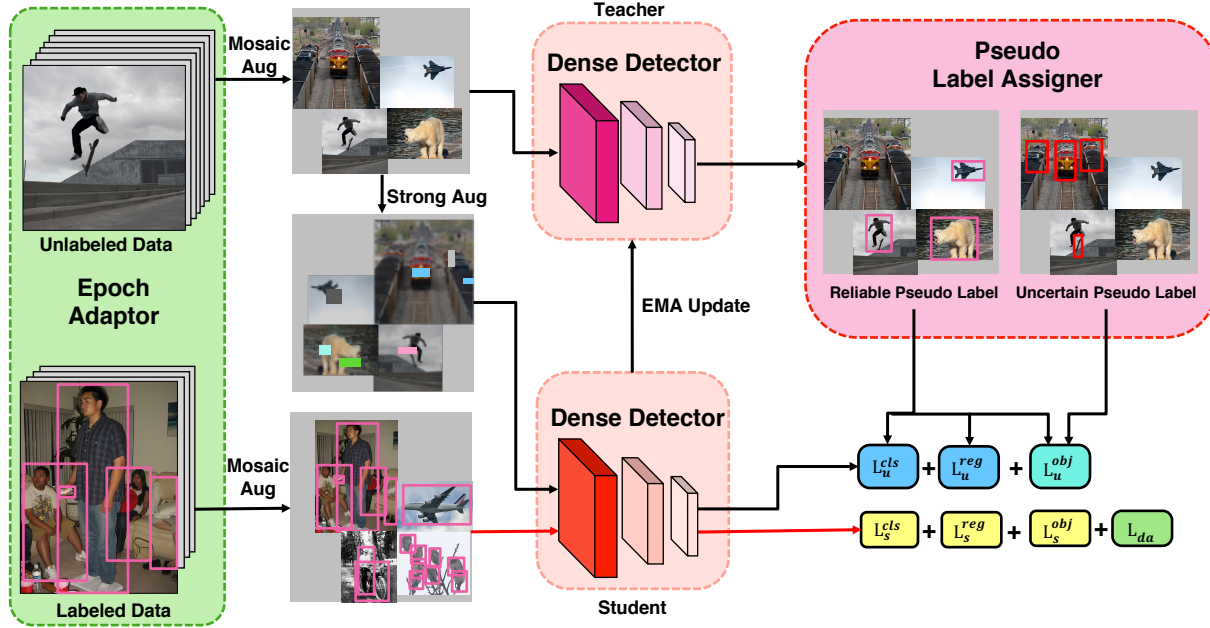


Figure 1. An overview of Efficient Teacher framework. Efficient Teacher proposes three modules to implement a scalable and effective SSOD framework, where Dense Detector improves the quality of pseudo labels with dense input while has better inference efficiency; Pseudo Label Assigner divides pseudo labels into two types to alleviate pseudo labels inconsistency problem; Epoch Adaptor reduces training time and the inconsistency of features.

mainly in pursuit of better accuracy, but usually sacrifice training efficiency. Moreover, most previous works only focus on specific detector architecture, but the variety of real-world applications require faster iterative detector design with lower compute resource and higher accuracy.

In this paper, what we pursue is to design a scalable and effective SSOD framework on an one-stage anchor-based detector while considering both inference and training efficiency. We add the effective techniques used in the YOLO series [2, 15, 42] to a classical RetinaNet [25] to design a new representative one-stage anchor-based detector baseline, called Dense Detector. We attempt to transplant a mature SSOD scheme, the Unbiased Teacher [27], to Dense Detector but find only 1.65  $AP_{50:95}$  improvement compared to the supervised method (shown in Table 5), which confirms the second problem mentioned above. According to design paradigm of the Dense Detector, we propose the Efficient Teacher framework to overcome these challenges in SSOD. Pseudo Label Assigner (PLA) is introduced to alleviate pseudo label inconsistency by exploiting a fine-grained pseudo label assignment strategy on the objectness branch design. By distinguishing the pseudo labels into the reliable and the uncertain regions, different loss calculation methods are used respectively. In addition, we propose Epoch Adaptor (EA), an end-to-end training strategy for SSOD which consists of a Burn-In stage and a semi-supervised training stage. In the Burn-In stage, only labeled data is used

for model parameter updating warm-up to obtain a student model with basic pseudo labels generation capability. Specially, the student model performs adversarial learning in the Burn-In stage so that the feature map of output is consistent on both labeled and unlabeled data. In the semi-supervised training stage, EA automatically calculates low and high threshold used by PLA according to the number of the ground truth (GT) after Mosaic augmentation in an epoch, which leads to stable and effective training. The main contributions of this paper are as follows:

- We design Dense Detector as a baseline model to compare the differences between YOLOv5 and RetinaNet, which leads to a performance improvement of 5.36  $AP_{50:95}$  by utilizing dense sampling.
- We propose an effective SSOD training framework called Efficient Teacher, which includes a novel pseudo label assignment mechanism, Pseudo Label Assigner, reducing the inconsistency of pseudo labels, and Epoch Adaptor, enabling a fast and efficient end-to-end semi-supervised training schedule.
- our experiments demonstrate that utilizing Efficient Teacher on YOLOv5 produces state-of-the-art results on VOC, COCO-standard, and COCO-additional datasets while consuming significantly fewer FLOPs than previous approaches.

## 2. Related Work

**Semi-supervised Object Detection.** Semi-supervised object detection, inherited from the semi-supervised image classification methods [1, 33, 35, 39, 43], is divided into consistency-based schemes [18, 37] and pseudo-labeling schemes [27, 36, 45, 50]. The latter has become the current mainstream approach. STAC [36] exploits weak and strong data augmentation to process unlabeled data respectively. Unbiased Teacher [27] follows a student-teacher mutual learning to generate more accurate pseudo labels. To balance the effect of pseudo labels, Soft Teacher [45] uses the scores of the pseudo labels as the weights for loss calculation. DSL [5] is the first attempt to perform semi-supervised training on an anchor-free detector (FCOS) [40]. To relieve inconsistency problems, LabelMatch [4] utilizes label distribution to dynamically determine the filtering threshold of different categories of pseudo labels. The methods above have been proven great performance on two-stage and anchor-free detectors, but can not perform well on an one-stage anchor-based detector. Our Efficient Teacher is proposed to bridge the gap between semi-supervised training and one-stage anchor-based detectors.

**Label Assignment.** Label assignment is the key component that determines the performance of an object detector. Many works have been proposed to improve the label assignment mechanism, such as ATSS [46], PAA [21], AutoAssign [51] and OTA [14]. Some researches [4] [28] have noticed that the default label assignment mechanism using in supervised object detection can not be applied in SSOD directly, which results in performance degradation. In this paper, we propose a novel pseudo label assignment that can adapt to SSOD training for one-stage anchor-based detectors.

**Domain Adaptation in Object Detection.** The task of domain-adaptive object detection [7, 10, 24, 41], aims to address the problem of domain shift [8]. The work [13] utilizes adversarial learning by training a discriminator with a gradient reverse layer (GRL) to generate domain-invariant feature. The work [10] introduces semi-supervised techniques used in Mean Teacher to alleviate domain bias, which reveals that domain shift is intrinsically related to inconsistency of semi-supervised task. This inspires Efficient Teacher to introduce adversarial learning in domain adaptation to alleviate the pseudo label inconsistency of SSOD training.

## 3. Efficient Teacher

Efficient Teacher is a novel and efficient framework for semi-supervised object detection, which significantly enhances the performance of one-stage anchor-based detectors. The framework is based on a student-teacher mutual learning approach, as shown in Figure 1, inspired by

Method	Resolution	Mosaic	Param.	FLOPs	$AP_{50:95}(\%)$
Faster R-CNN [31]	[1333,800]		39.8M	202.31G	40.3
FCOS [40]	[1333,800]		32.02M	200.59G	38.5
YOLOv5 <i>w/o</i>	[640,640]		46.56M	109.59G	41.2
YOLOv5 [19]	[640,640]	✓	46.56M	109.59G	47.87
YOLOv7 [42]	[640,640]	✓	37.62M	106.59G	51.5
RetinaNet [25]	[1333,800]		37.74M	239.32G	39.5
Dense Detector	[640,640]	✓	42.13M	169.61G	44.86

Table 1. Comparison with Faster R-CNN, FCOS, YOLOv5, YOLOv7, RetinaNet and Dense Detector. The top section shows results for object detectors without Mosaic augmentation, the middle section shows results with Mosaic augmentation during training. Dense Detector achieves comparable results to RetinaNet baseline, having lower FLOPs but greatly improved  $AP_{50:95}$ . Both Faster R-CNN, FCOS, RetinaNet and Dense Detector uses ResNet-50-FPN as backbone.  $AP_{50:95}$  is reported on COCO val dataset.

previous works [4, 5, 27, 45]. Our proposed Pseudo Label Assigner method divides pseudo labels into reliable and uncertain ones based on their scores, with reliable pseudo labels used for default supervised training, and uncertain ones used to guide the training of the student model with objectness scores. The Epoch Adaptor method is used to speed up convergence by performing domain adaptation between labeled and unlabeled data, and switching the main epoch from labeled to unlabeled data after a burn-in stage. Throughout the training process, the teacher model employs the Exponential Moving Average (EMA) technique for updates.

### 3.1. Dense Detector

The imbalance of positive and negative training assignment samples often leads to unqualified pseudo labels in one-stage anchor-based detectors. To address this issue, we introduce the design of dense inputs. YOLOv5 [19] is a widely-used one-stage anchor-based detector in industry due to its friendly-deployed support and fast training speed. Results in Table 1 demonstrate that YOLOv5 *w/o* outperforms RetinaNet in terms of performance and computation. Furthermore, with dense image inputs after Mosaic augmentation, the  $AP_{50:95}$  of YOLOv5 can be boosted from 41.2 to 47.87. YOLOv7 further improves the  $AP_{50:95}$  to 51.5 with the help of dense flow of information and gradients on the basis of dense inputs. Our study suggests that improvements in the performance of one-stage anchor-based detectors often require dense inputs. Therefore, we propose Dense Detector as a base model for SSOD under dense inputs.

Dense Detector is modified from RetinaNet with ResNet-50-FPN backbone while changing the number of FPN output from 5 to 3, eliminating the weight sharing between detection headers and reducing the input resolution from 1333 to 640 for both training and inference. What's

more, Dense Detector has two outputs(a classification score and a bounding-box offset) and a third branch that outputs the objectness score. Dense Detetor gained 5.36%  $AP_{50:95}$  boost and 30% faster relative to RetinaNet as shown in Table 1. Specifically, Dense Detector obtains objectness score by calculating the Complete Intersection over Union(CIoU) [48] between the predicted and GT boxes. Objectness output directly indicates position response, and reflects the location quality of the predicted boxes. As illustrated in Figure 1, the pseudo labels in SSOD are the predicted boxes of unlabeled data, the objectness scores of which indicate the location quality of pseudo labels. Thus, compared to RetinaNet with only a classification branch, Dense Detector with an extra objectness branch can indicate the location quality of pseudo labels during SSOD training as shown in Figure 2.

To verify the performance of Dense Detector in SSOD, we apply the classic SSOD method(Unbiased Teacher [27]) to the Dense Detector, which contains labeled and unlabeled data, teacher and student model, and a pseudo label filter to select pseudo labels. Furthermore, both labeled and unlabeled data branches adopt loss definition in Equation 2. However, in contrast to Unbiased Teacher on Faster R-CNN in Table 2, the  $AP_{50:95}$  improvement of Unbiased Teacher drops from 7.64 to 4.3 on Dense Detector. This motivated us to develop the following Pseudo Label Assigner that plays a key role in pseudo label assignment.

### 3.2. Pseudo Label Assigner

The core problem in SSOD is how to assign pseudo labels, as sub-optimal assignments can lead to inconsistent pseudo labels and deteriorating performance of the mutual learning mechanism. Pseudo Label Filter (Figure 3) is a naive approach that assigns pseudo labels by setting a threshold, with those below the threshold being considered as background. However, this method can result in sub-optimal assignments, as shown in Figure 3: in the top case, the pseudo label with red color is incorrectly assigned due to a threshold value of 0.1, leading to the student learning the wrong label; in the bottom case, the threshold of 0.6 excludes low-scoring pseudo labels but incorrectly treats the pseudo label in red color as background, resulting in decreased training for this class. These cases highlight the need for more reasonable strategies for pseudo label assignment in Dense Detector.

Proposed in this work, Pseudo Label Assigner (PLA) provides a more refined assignment of the pseudo labels generated by Dense Detector. In PLA, pseudo labels obtained after Non-Maximum Suppression (NMS) are separated into two categories: reliable and uncertain pseudo labels. The high and low threshold  $\tau_1, \tau_2$  of the pseudo label score is used to determine two types of pseudo labels. Pseudo labels with scores between  $\tau_1, \tau_2$  are considered un-

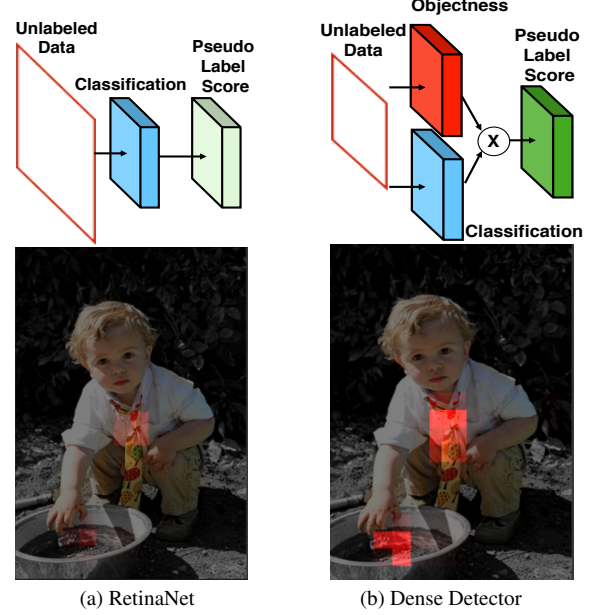


Figure 2. Comparison of pseudo label score heatmaps from RetinaNet and Dense Detector. Darker color indicates higher score. (a) RetinaNet produces sparse response due to the calculation of classification scores from pseudo labels generated from unlabeled data of  $1333 \times 800$  input resolution. (b) Dense Detector, with input resolution of  $640 \times 640$ , uses a weighted pseudo label score based on objectness and classification scores, resulting in a more robust and dense response..

certain, and ignoring the loss of these labels directly results in improved performance on Dense Detector, as shown in Table 5. In addition to solving the sub-optimal problem caused by Pseudo Label Filter, PLA includes an unsupervised loss that efficiently leverages uncertain pseudo labels. The loss of Dense Detector in SSOD is defined as a pair of single labeled image and single unlabeled image:

$$L = L_s + \lambda L_u \quad (1)$$

where  $L_s$  represents the loss function computed on a labeled image, while  $L_u$  represents the loss function computed on an unlabeled image,  $\lambda$  is used to balance the supervised loss and the semi-supervised loss, which is set to 3.0 in this paper. The  $L_s$  follows the standard loss function in [19]:

$$L_s = \sum_{h,w} (CE(X_{(h,w)}^{cls}, Y_{(h,w)}^{cls}) + CIoU(X_{(h,w)}^{reg}, Y_{(h,w)}^{reg}) + CE(X_{(h,w)}^{obj}, Y_{(h,w)}^{obj})) \quad (2)$$

where CE indicates cross-entropy loss function,  $X_{(h,w)}$  is the output of student model, and  $Y_{(h,w)}$  means the sampled results generated by the label assigner of Dense Detector. The  $L_u$  is defined as follows:

$$L_u = L_u^{cls} + L_u^{reg} + L_u^{obj} \quad (3)$$

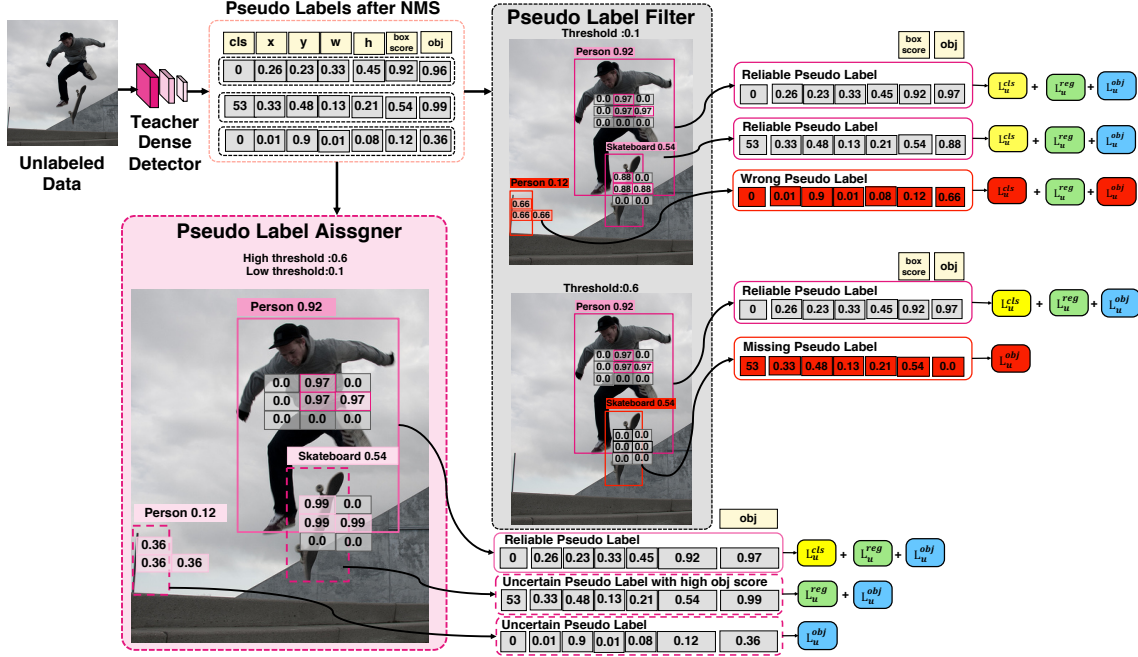


Figure 3. Comparison of the impact of different pseudo label select strategy. In Pseudo Label Filter, a widely-used method in SSOD [27,45], Setting the threshold too low (0.1) can result in the generation of incorrect pseudo labels (indicated by the red line), while a threshold that is too high (0.6) may exclude reliable pseudo labels. This can lead to suboptimal assignments and adversely affect the training of the network. To address this issue, we propose the Pseudo Label Assigner method, which categorizes pseudo labels into reliable and uncertain categories based on high and low thresholds, respectively. The uncertain pseudo labels are assigned soft labels as targets for  $L_u^{obj}$  to improve the quality of pseudo labels in SSOD.

$$L_u^{cls} = \sum_{h,w} (\mathbb{1}_{\{p(h,w) \geq \tau_2\}} CE(X_{(h,w)}^{cls}, \hat{Y}_{(h,w)}^{cls})) \quad (4)$$

$$L_u^{reg} = \sum_{h,w} (\mathbb{1}_{\{p(h,w) \geq \tau_2 \text{ or } obj(h,w) > 0.99\}} CIOU(X_{(h,w)}^{reg}, \hat{Y}_{(h,w)}^{reg})) \quad (5)$$

$$L_u^{obj} = \sum_{h,w} (\mathbb{1}_{\{p(h,w) \leq \tau_1\}} CE(X_{(h,w)}^{obj}, \mathbf{0}) + \mathbb{1}_{\{p(h,w) \geq \tau_2\}} CE(X_{(h,w)}^{obj}, \hat{Y}_{(h,w)}^{obj}) + \mathbb{1}_{\{\tau_1 < p(h,w) < \tau_2\}} CE(X_{(h,w)}^{obj}, \hat{obj}_{(h,w)})) \quad (6)$$

where  $\hat{Y}_{(h,w)}^{cls}$ ,  $\hat{Y}_{(h,w)}^{reg}$ ,  $\hat{Y}_{(h,w)}^{obj}$  is the classification score, regression, objectness score of sampled results from PLA at location  $(h, w)$  on feature map separately.  $\hat{obj}_{(h,w)}$  is the objectness score of pseudo label at  $(h, w)$ .  $p(h, w)$  is the score of pseudo label at  $(h, w)$ .  $\mathbb{1}_{\{\cdot\}}$  is the indicator function, which outputs 1 if condition  $\{\cdot\}$  is satisfied and 0 otherwise.

PLA categorizes uncertain pseudo labels into two types: 1) those with high classification scores and 2) those with

high objectness scores, which are handled differently. For the first type, only  $L_u^{obj}$  is calculated, and the targets of the cross-entropy  $\hat{Y}_{(h,w)}^{obj}$  are replaced with the soft label  $\hat{obj}_{(h,w)}$ , indicating that these pseudo labels are not classified as either background or positive samples. For the second type, PLA calculates  $L_u^{reg}$  when the objectness score is greater than 0.99, as these pseudo labels have good regression results but insufficient classification scores to determine their label category. PLA aims to convert more uncertain pseudo labels into true positives using  $L_u^{reg}$ , as more than 70% of uncertain pseudo labels are false positives due to inaccurate prediction boxes during SSOD training on COCO. Therefore, PLA suppresses the inconsistency of pseudo labels through a soft label learning mechanism, without affecting the loss of reliable pseudo labels. Further details can be found in the Appendix.

### 3.3. Epoch Adaptor

Although the PLA addresses the issue of pseudo label inconsistency in SSOD, the hyperparameters  $\tau_1$  and  $\tau_2$  are still hand-tuned and influenced by the distribution of labeled data. Additionally, the training scheme significantly impacts the performance of SSOD method. Current SSOD methods use two types of training scheme: alternating train-

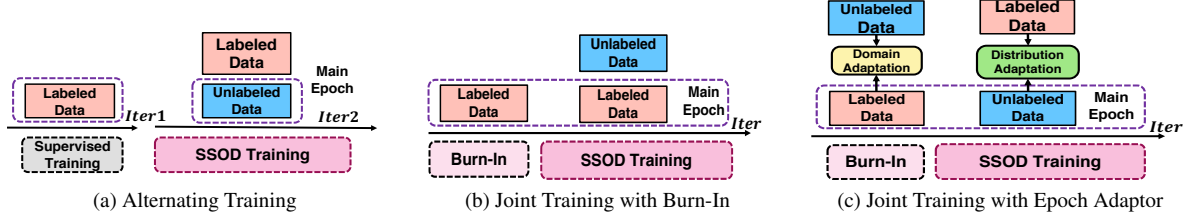


Figure 4. Training strategies for Dense Detector: (a) supervised training on labeled data followed by SSOD training on unlabeled data; (b) supervised training on labeled data with additional SSOD training on unlabeled data; (c) end-to-end training on both labeled and unlabeled data with Epoch Adaptor incorporating Domain and Distribution Adaptation for improved convergence and feature distribution.

ing and joint training with Burn-In. Alternating training is time-consuming, but it allows for the precise calculation of the reliable pseudo label threshold [4]. In contrast, joint training with Burn-In directly adds unlabeled data to the training process and has a faster training speed but reduces SSOD training stability [27, 45]. Soft Teacher [45] suggested artificially increasing the ratio of labeled and unlabeled data in a mini-batch (1:4) due to the differences in the number of labeled and unlabeled data in the training set.

Building upon the joint training scheme, we present the Epoch Adaptor(EA) framework that enhances the Burn-In stage’s training efficiency and consistency. Moreover, EA dynamically calculate the  $\tau_1$  and  $\tau_2$  parameters of PLA in the SSOD training stage. During the Burn-In stage, EA utilizes domain adaptation to allow Dense Detector to train simultaneously on labeled and unlabeled data. Specifically, EA incorporates a 1:1 ratio of labeled and unlabeled data in 1:1 ratio and trains a domain adaptation classifier on each pair of labeled and unlabeled data. The domain adaptation loss function as follow:

$$L_{da} = - \sum_{h,w} [D \log p(h,w) + (1-D) \log(1-p(h,w))]. \quad (7)$$

where  $p(h,w)$  is the output of the domain classifier.  $D = 0$  for labeled data and  $D = 1$  for unlabeled data. We use the gradient reverse layer (GRL) [13], whereas the ordinary gradient descent is applied for training the domain classifier and the sign of the gradient is reversed when passing through the GRL layer to optimize the base network. In Burn-In stage, the supervised loss in one image can be reformulated as follows:

$$L_s = \sum_{h,w} (CE(X_{(h,w)}^{cls}, Y_{(h,w)}^{cls}) + CIOU(X_{(h,w)}^{reg}, Y_{(h,w)}^{reg}) + CE(X_{(h,w)}^{obj}, Y_{(h,w)}^{obj})) + \lambda L_{da} \quad (8)$$

where  $\lambda$  is the hyper-parameter to control the contribution of domain adaptation, which is 0.1 in our experiments. The expression capability of the model is enhanced by allowing the detector to see the unlabeled data in Burn-In.

During the SSOD training stage, EA applies a strategy where the current main epoch switches from the labeled data to the unlabeled data, as depicted in Figure

4. To dynamically calculate the  $\tau_1$  and  $\tau_2$  thresholds of PLA in each epoch, we implement a distribution adaptation method based on the re-distribution method in Label-Match [4]. This is necessary because Mosaic augmentation increases the number of ground truth (GT) annotations in labeled data, which alters the GT counting setting used in the re-distribution method based on offline annotations in LabelMatch [4]. Specifically, the GT count per image increases from 7.24 to 21.4 during SSOD training on the COCO dataset. The  $\tau_1$  and  $\tau_2$  thresholds at the  $k$ -th epoch are determined as follows:

$$\tau_1^k = P_c^k [n_c^k \cdot \frac{N_u}{N_l}] \quad (9)$$

$$\tau_2^k = P_c^k [\alpha \% \cdot n_c^k \cdot \frac{N_u}{N_l}], \quad (10)$$

The reliable ratio  $\alpha$  is set to 60 for all experiments, and  $P_c^k$  represents the list of pseudo label scores of the  $c$ -th class at the  $k$ -th epoch. Meanwhile,  $N_l$  and  $N_u$  denote the number of labeled and unlabeled data, and  $n_c^k$  represents the number of  $c$ -th class ground truth annotations that are counted by EA at the  $k$ -th epoch. By dynamically calculating the appropriate thresholds at each epoch, EA enables joint training to be more adaptable to dynamic data distributions. This approach differs from the offline statistical method used by LabelMatch [4].

By enabling the SSOD training to obtain higher  $AP_{50:95}$  and faster convergence without the need for hand-tuned thresholds, EA effectively reduces the overall training time of joint training scheme. The experimental results demonstrating these effects are presented in Section 4.

## 4. Experiments

### 4.1. Experimental Setup

**Datasets.** We validate our method on MS-COCO [26] and VOC [12] benchmarks: (1) COCO-standard: 1%, 2%, 5%, 10% of the images are sampled on COCO as labeled data, and all the remaining data are used as unlabeled data. (2) COCO-additional: train2017 dataset is set as the labeled dataset and COCO2017-unlabeled is as the unlabeled dataset. (3) VOC: VOC07 trainval data is as the labeled

Method		%1	%2	%5	%10	FLOPs
Two-stage anchor-based	Supervised	9.05	12.70	18.47	23.86	202.31G
	STAC [36]	13.97 $\pm$ 0.35(+4.92)	18.25 $\pm$ 0.25 (+5.91)	24.38 $\pm$ 0.12 (+5.91)	28.64 $\pm$ 0.21 (+4.78)	202.31G
	Instant Teaching [50]	18.05 $\pm$ 0.15 (+9.00)	22.45 $\pm$ 0.15 (+9.75)	26.75 $\pm$ 0.05 (+8.28)	30.40 $\pm$ 0.05 (+6.54)	202.31G
	Humber teacher [38]	16.96 $\pm$ 0.38 (+7.91)	21.72 $\pm$ 0.24 (+9.02)	27.70 $\pm$ 0.15 (+9.23)	31.61 $\pm$ 0.28 (+7.75)	202.31G
	Unbiased Teacher [27]	20.75 $\pm$ 0.12 (+11.70)	24.30 $\pm$ 0.07 (+9.80)	28.27 $\pm$ 0.11 (+9.80)	31.50 $\pm$ 0.10 (+7.64)	204.13G
	Soft Teacher [45]	20.46 $\pm$ 0.39 (+11.41)	-	30.74 $\pm$ 0.08 (+12.27)	34.04 $\pm$ 0.14 (+10.18)	202.31G
	LabelMatch [4]	<b>25.81</b> $\pm$ 0.28 (+16.76)	-	32.70 $\pm$ 0.18 (+14.23)	35.49 $\pm$ 0.17 (+11.63)	202.31G
	PseCo [23]	22.43 $\pm$ 0.36 (+13.38)	27.77 $\pm$ 0.18 (+15.07)	32.50 $\pm$ 0.08 (+14.03)	36.06 $\pm$ 0.24 (+12.20)	202.31G
One-stage anchor-free	Supervised	9.53	11.71	18.74	23.70	200.59G
	Unbiased Teacher v2 [28]	22.71 $\pm$ 0.42 (+13.18)	26.03 $\pm$ 0.12 (+14.32)	30.08 $\pm$ 0.04 (+11.34)	32.61 $\pm$ 0.03 (+8.91)	200.59G
	DSL [5]	22.03 $\pm$ 0.28 (+12.50)	25.19 $\pm$ 0.37 (+13.48)	30.87 $\pm$ 0.24 (+12.13)	36.22 $\pm$ 0.18 (+12.52)	200.59G
	Dense Teacher [49]	22.38 $\pm$ 0.31 (+12.85)	27.20 $\pm$ 0.20 (+15.49)	33.01 $\pm$ 0.21 (+14.27)	37.13 $\pm$ 0.12 (+13.43)	200.59G
One-stage anchor-based	Supervised	11.29	13.12	20.28	26.04	169.61G
	Unbiased Teacher* [27]	18.81 $\pm$ 0.28 (+7.52)	22.72 $\pm$ 0.21 (+9.60)	28.35 $\pm$ 0.12 (+8.15)	30.34 $\pm$ 0.09 (+4.30)	169.61G
	Ours	20.18 $\pm$ 0.21 (+8.87)	25.85 $\pm$ 0.13 (+12.73)	30.41 $\pm$ 0.08 (+10.13)	33.44 $\pm$ 0.11 (+7.40)	169.61G
	Ours $\dagger$	23.76 $\pm$ 0.13 (+12.47)	<b>28.70</b> $\pm$ 0.14 (+15.58)	<b>34.11</b> $\pm$ 0.09 (+13.83)	<b>37.90</b> $\pm$ 0.04 (+11.86)	109.59G

Table 2. Experimental results on COCO-standard ( $AP_{50:95}$ ), \* means re-implemented results on Dense Detector,  $\dagger$  means Efficient Teacher with YOLOv5l [19]. All the results are the average of 5 folds.

dataset and VOC12 trainval is used as the unlabeled dataset. We adopt the mean average precision  $AP_{50:90}$  as the evaluation metric.

**Network.** To verify that our proposed method is scalable, we used three Dense Detector architectures: The first one uses ResNet-50-FPN in Dense Detector. The second one replaces the original backbone with CSPNet and the Neck with PAN, which is similar with YOLOv5. The last one follows YOLOv7 to add ELAN module and RepConv [11] into backbone.

**Implementation Details.** We use 8 NVIDIA-V100 GPUs with 16G memory per GPU. We randomly sample 32 images from labeled data and 32 images from unlabeled data with ratio 1:1 in each iteration. For training configurations, the learning rate is 0.01 all the time, the  $\tau_1$  and  $\tau_2$  are calculated by EA. We used both weak and strong data augmentation. Mosaic is used in weak data augmentation. In the strong data augmentation, Mosaic, left-right flip, large scale jittering, graying, Gaussian blur, cutout, and color space conversion are selected. The max epoch is 300. Smoothing hyper-parameter in EMA is 0.999.

## 4.2. Results

**COCO-standard.** In Table 2, we validate our proposed method on COCO-standard and the performance of Efficient Teacher is better than Unbiased Teacher on Dense Detector. Moreover, when applying standard YOLOv5l Backbone, our method achieves state-of-the-art results on labeled data with 2%, 5%, 10% coefficients. Comparing to previous state-of-the-arts, Efficient Teacher is the second highest in 1% labeled setting, after LabelMatch [4], but greatly improved both in 5% and 10% COCO using fewer FLOPs.

**COCO-additional.** Results in Table 3 show our proposed method on COCO-additional, the gain effect of Effi-

cient Teacher is better than Unbiased Teacher, which shows 1.01 increase on  $AP_{50:95}$ . Backbone in all experiments is YOLOv5l.

Method	$AP_{50:95}$
Supervised $\dagger$	47.87
Unbiased Teacher [27] $\dagger$	48.48(+0.61)
Ours $\dagger$	<b>48.88(+1.01)</b>

Table 3. Experimental results on COCO-additional.

**PASCAL-VOC.** Table 4 shows the results of experiments conducted on VOC are convincing. Our method achieves 58.30 on  $AP_{50:95}$ , which shows more accurate bounding box regression results.

Method	$AP_{50:95}$	$AP_{50}$
STAC [36]	44.64	77.45
Instant Teacher [50]	50.00	79.20
Unbiased Teacher [27]	48.69	77.37
Dense Teacher [49]	55.87	79.89
DSL [5]	56.80	80.70
Unbiased Teacher v2 [28]	56.87	81.29
Ours $\dagger$	<b>58.30</b>	<b>81.60</b>

Table 4. Experimental results on PASCAL-VOC.

## 4.3. Ablation Studies

In ablation studies, we conducted experiments using the 10% COCO-standard dataset (one of 5 folds). The backbone is YOLOv5l.

**Effect of Pseudo Label Assigner.** The impact of the proposed Pseudo Label Assigner is presented in Table 5. We observe that applying the Unbiased Teacher method to the Dense detector with a threshold of 0.3 for pseudo label generation only leads to a modest  $AP_{50:95}$  improvement of 1.65, which is considerably lower than the 7.6  $AP_{50:95}$  gain achieved by the Unbiased Teacher [27] applied to the Faster R-CNN. When neglecting the uncertain pseudo labels, the  $AP_{50:95}$  further increases to 35.2. However, by utilizing the Pseudo Label Assigner to handle the uncertain pseudo labels, we obtain a significant improvement of 7.45 in  $AP_{50:95}$ , resulting in a final performance of 37.90, which is comparable to that of the Unbiased Teacher applied to the Faster R-CNN.

Method	$AP_{50:95}$	$AP_{50}$
Supervised	30.45	44.65
Unbiased Teacher [27]	32.10 (+1.65)	47.30 (+2.65)
Ignore uncertain pseudo label [5]	35.20 (+4.75)	52.00 (+7.35)
Pseudo Label Assigner	<b>37.90 (+7.45)</b>	<b>54.19 (+9.54)</b>

Table 5. Ablation study about different pseudo label assignment methods.

**Effect of dynamic threshold.** We evaluate the impact of varying the threshold value  $\tau_2$  in the Pseudo Label Assigner method on the COCO 10% standard task. As shown in Table 6, increasing the value of  $\tau_2$  results in a decreasing trend for  $AP_{50:95}$ , which is indicative of fewer reliable pseudo labels and more uncertain ones. This highlights the importance of having an appropriate balance between reliable and uncertain pseudo labels, as the decrease in the number of reliable pseudo labels can negatively impact the effectiveness of SSOD training. Notably, the use of Epoch Adaptor to dynamically calculate the value of  $\tau_2$  results in the best performance without requiring manual tuning efforts.

$\tau_2$	$AP_{50:95}$	$AP_{50}$
0.4	37.20	54.08
0.5	37.20	54.10
0.6	36.90	53.77
0.7	35.10	51.60
EA	<b>37.90</b>	<b>54.80</b>

Table 6. Ablation studies on dynamic threshold, EA indicates  $\tau_2$  is calculated by Epoch Adaptor.

**Effect of the Epoch Adaptor.** We conducted an ablation study of Epoch Adaptor, as shown in Figure 5. The results demonstrate that after jointly training with Epoch Adaptor, the network achieves superior performance with fewer iterations compared to the other two SSOD train-

ing methods, namely fully supervised and alternating training, which involves training a detector for 70K iterations in a fully supervised manner followed by 240K iterations of semi-supervised training.

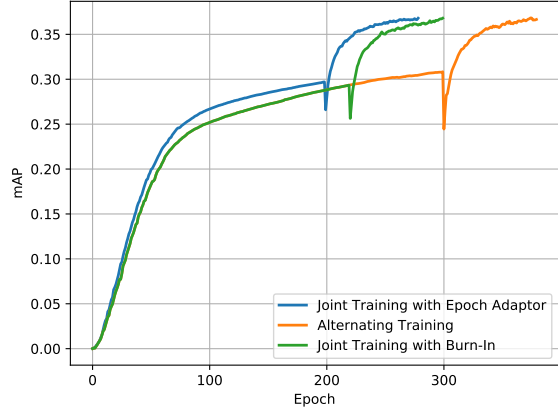


Figure 5. Performance ( $AP_{50:95}$ ) comparisons of Epoch Adaptor, Alternating Training and Joint Training with Burn-In methods on COCO standard 10%.

## 5. Conclusion

In this paper, we present Efficient Teacher, a method to bridge the gap between SSOD and one-stage anchor-based detectors, by building on the efficient dense input handling of Dense Detector. Our approach introduces the Pseudo Label Assigner to effectively utilize both reliable and uncertain pseudo labels, based on an analysis of their assignment in SSOD. In addition, we introduce Epoch Adaptor, a training scheme that maximizes the efficiency of training and utilization of both labeled and unlabeled data. Our approach achieves state-of-the-art results on VOC, COCO-standard, and COCO-additional datasets.

## References

- [1] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A Raffel. Mixmatch: A holistic approach to semi-supervised learning. *Advances in neural information processing systems*, 32, 2019. 1, 3
- [2] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934*, 2020. 1, 2
- [3] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6154–6162, 2018. 1
- [4] Binbin Chen, Weijie Chen, Shicai Yang, Yunyi Xuan, Jie Song, Di Xie, Shiliang Pu, Mingli Song, and Yueting

- Zhuang. Label matching semi-supervised object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14381–14390, 2022. 1, 3, 6, 7, 11, 12
- [5] Binghui Chen, Pengyu Li, Xiang Chen, Biao Wang, Lei Zhang, and Xian-Sheng Hua. Dense learning based semi-supervised object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4815–4824, 2022. 1, 3, 7, 8
- [6] Chaoqi Chen, Zebiao Zheng, Xinghao Ding, Yue Huang, and Qi Dou. Harmonizing transferability and discriminability for adapting object detectors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8869–8878, 2020. 12
- [7] 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. *arXiv preprint arXiv:2102.11614*, 2021. 3
- [8] Yuhua Chen, Wen Li, Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Domain adaptive faster r-cnn for object detection in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3339–3348, 2018. 3
- [9] 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. 12
- [10] 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. 3, 12
- [11] Xiaohan Ding, Xiangyu Zhang, Ningning Ma, Jungong Han, Guiguang Ding, and Jian Sun. Repvgg: Making vgg-style convnets great again. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13733–13742, 2021. 7
- [12] Mark Everingham and John Winn. The pascal visual object classes challenge 2012 (voc2012) development kit. *Pattern Anal. Stat. Model. Comput. Learn., Tech. Rep.*, 2007:1–45, 2012. 6
- [13] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning*, pages 1180–1189. PMLR, 2015. 3, 6
- [14] 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
- [15] Zheng Ge, Songtao Liu, Feng Wang, Zeming Li, and Jian Sun. Yolox: Exceeding yolo series in 2021. *arXiv preprint arXiv:2107.08430*, 2021. 1, 2
- [16] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3354–3361. IEEE, 2012. 12
- [17] Mengzhe He, Yali Wang, Jiaxi Wu, Yiru Wang, Hanqing Li, Bo Li, Weihao Gan, Wei Wu, and Yu Qiao. Cross domain object detection by target-perceived dual branch distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9570–9580, 2022. 12
- [18] Jisoo Jeong, Seungeui Lee, Jeessoo Kim, and Nojun Kwak. Consistency-based semi-supervised learning for object detection. *Advances in neural information processing systems*, 32, 2019. 3
- [19] G Jocher, A Chaurasia, A Stoken, J Borovec, NanoCode012, Y Kwon, TaoXie, J Fang, imyhxy, and K Michael. ultralytics/yolov5: v6. 1-tensorrt, tensorflow edge tpu and openvino export and inference. *Zenodo, Feb*, 22, 2022. 1, 3, 4, 7
- [20] 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? *arXiv preprint arXiv:1610.01983*, 2016. 12
- [21] Kang Kim and Hee Seok Lee. Probabilistic anchor assignment with iou prediction for object detection. In *European Conference on Computer Vision*, pages 355–371. Springer, 2020. 3
- [22] Chuyi Li, Lulu Li, Hongliang Jiang, Kaiheng Weng, Yifei Geng, Liang Li, Zaidan Ke, Qingyuan Li, Meng Cheng, Weiqiang Nie, et al. Yolov6: a single-stage object detection framework for industrial applications. *arXiv preprint arXiv:2209.02976*, 2022. 1
- [23] Gang Li, Xiang Li, Yujie Wang, Shanshan Zhang, Yichao Wu, and Ding Liang. Pseco: Pseudo labeling and consistency training for semi-supervised object detection. *arXiv preprint arXiv:2203.16317*, 2022. 1, 7
- [24] 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 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 8474–8481, 2021. 3
- [25] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 1, 2, 3
- [26] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 6
- [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. *arXiv preprint arXiv:2102.09480*, 2021. 1, 2, 3, 4, 5, 6, 7, 8, 11
- [28] Yen-Cheng Liu, Chih-Yao Ma, and Zsolt Kira. Unbiased teacher v2: Semi-supervised object detection for anchor-free and anchor-based detectors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9819–9828, 2022. 3, 7
- [29] 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. 12
- [30] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018. 1
- [31] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015. 1, 3
- [32] Kuniaki Saito, Yoshitaka Ushiku, Tatsuya Harada, and Kate Saenko. Strong-weak distribution alignment for adaptive object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6956–6965, 2019. 12
- [33] Mehdi Sajjadi, Mehran Javanmardi, and Tolga Tasdizen. Regularization with stochastic transformations and perturbations for deep semi-supervised learning. *Advances in neural information processing systems*, 29, 2016. 3
- [34] 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. 12
- [35] Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel, Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in neural information processing systems*, 33:596–608, 2020. 1, 3
- [36] 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. 3, 7
- [37] Peng Tang, Chetan Ramaiah, Yan Wang, Ran Xu, and Caiming Xiong. Proposal learning for semi-supervised object detection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2291–2301, 2021. 3
- [38] Yihe Tang, Weifeng Chen, Yijun Luo, and Yuting Zhang. Humble teachers teach better students for semi-supervised object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3132–3141, 2021. 7
- [39] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in neural information processing systems*, 30, 2017. 3
- [40] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9627–9636, 2019. 1, 3
- [41] 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. 3, 12
- [42] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors. *arXiv preprint arXiv:2207.02696*, 2022. 1, 2, 3
- [43] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student improves imagenet classification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10687–10698, 2020. 3
- [44] Minghao Xu, Hang Wang, Bingbing Ni, Qi Tian, and Wenjun Zhang. Cross-domain detection via graph-induced prototype alignment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12355–12364, 2020. 12
- [45] 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. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3060–3069, 2021. 1, 3, 5, 6, 7
- [46] Shifeng Zhang, Cheng Chi, Yongqiang Yao, Zhen Lei, and Stan Z Li. Bridging the gap between anchor-based and anchor-free detection via adaptive training sample selection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9759–9768, 2020. 1, 3
- [47] Yueming Zhang, Xingxu Yao, Chao Liu, Feng Chen, Xiaolin Song, Tengfei Xing, Runbo Hu, Hua Chai, Pengfei Xu, and Guoshan Zhang. S4od: Semi-supervised learning for single-stage object detection. *arXiv preprint arXiv:2204.04492*, 2022. 1
- [48] Zhaohui Zheng, Ping Wang, Dongwei Ren, Wei Liu, Rongguang Ye, Qinghua Hu, and Wangmeng Zuo. Enhancing geometric factors in model learning and inference for object detection and instance segmentation. *IEEE Transactions on Cybernetics*, 2021. 4
- [49] Hongyu Zhou, Zheng Ge, Songtao Liu, Weixin Mao, Zeming Li, Haiyan Yu, and Jian Sun. Dense teacher: Dense pseudo-labels for semi-supervised object detection. *arXiv preprint arXiv:2207.02541*, 2022. 1, 7
- [50] 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, 3, 7
- [51] 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. More Analysis of Efficient Teacher

In this section, we aim to identify real gain of Efficient Teacher and make an in-depth study of it. Efficient Teacher is conceptually simple: it uses Dense Detector to generate high quality pseudo labels, Pseudo Label Assigner to alleviate pseudo label inconsistency problem and Epoch Adaptor to count the number of labels per image.

The backbone of Dense Detector in Efficient Teacher adopts CSPNet and PAN and all experiments are on one of five folds COCO standard 10% dataset.

### A.1. Analysis of Dense Detector

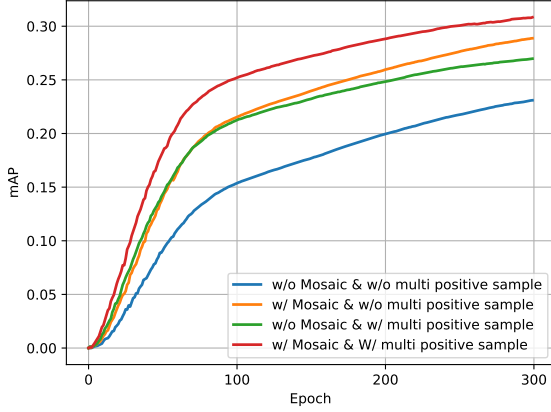


Figure 6. comparison between performance of different input and sample strategies, the result shows stable rise with dense inputs.

As we explain in Section 3.1, the dense inputs of Dense Detector is an engineering technique to balances the importance of positive/negative examples. The dense inputs consists of two parts: dense image inputs and dense sample inputs, we design the ablation experiments about dense image inputs(Mosaic augmentation) and dense sample inputs(multi positive sample) in Dense Detector. The Mosaic augmentation means patch four images into one and multi positive sample extend label assignment sample to adjacent points on feature map. We report the supervised training performance of Dense Detector, in which the higher mAP corresponds the better pseudo label quality in SSOD. It is revealed in Figure 6 that dense image inputs(Mosaic augmentation) shows massive improvement on mAP from 23.1 to 28.9 and dense sample inputs(multi positive sample) improves the mAP from 23.1 to 26.9, further more, with both two techniques the mAP can be boosted from 23.1 to 30.5. The dense inputs of Dense Detector has a positive effect on the efficiency and accuracy of Efficient Teacher.

### A.2. Analysis of Pseudo Labels Assigner

The Pseudo Label Assigner divides pseudo labels into two types: reliable and uncertain ones. In this part, we present more analysis about ratio of true positive and false positive in pseudo labels. Figure 7 shows the average statistics of reliable and uncertain pseudo labels in one epoch, the True Positive means the pseudo labels have the same class as ground truth and an IoU overlap greater than 0.5, the Loc False Positive means IoU overlap with ground truth less than 0.5 and Cls False Positive indicates pseudo labels is misclassified. What can be clearly seen in figures is 76% of reliable pseudo labels are True Positive while only 24%

of uncertain pseudo labels comply with the requirements. The misclassified pseudo labels are only small percentage of both reliable and uncertain pseudo labels, meanwhile, the weight of poor location in uncertain pseudo labels reaches 70% which motivated us to design  $L_u^{reg}$  in Section 3.2.

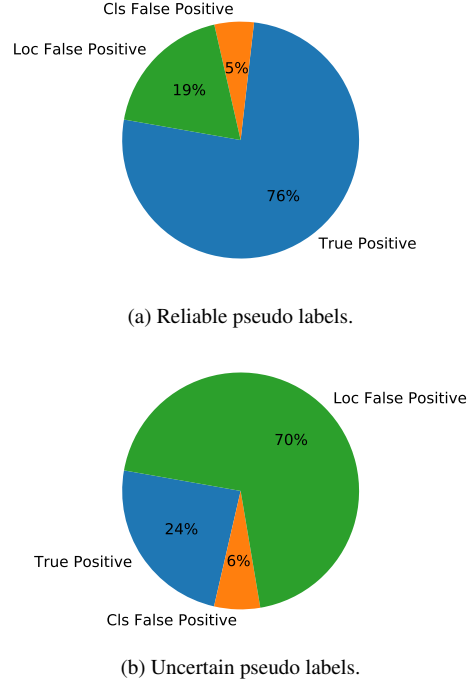


Figure 7. Weight of true and false positive compared across reliable and uncertain pseudo labels in Efficient Teacher.

### A.3. Analysis of Epoch Adaptor

We propose Epoch Adaptor to calculate the distribution of ground truth in each epoch due to Mosaic augmentation. We perform the qualitative comparisons between ground truth distribution of the proposed Efficient Teacher and Unbiased Teacher [27] method. The contents of images and instance distribution are shown in Figure 8, in which Unbiased Teacher has monotonous and large instance while Efficient Teacher has various and small instance. Moreover, the Mosaic augmentation is executed online which disabled the offline label number calculation in LabelMatch [4]. As discussed in Section 3.3, we observe the pseudo label and ground truth in Efficient Teacher and implement a distribution adaptation method to dynamically calculate number of ground truth per image which determines the  $\tau_1$  and  $\tau_2$  in Pseudo Label Assigner. Details are shown in Figure 9.

## B. Experiments on CityScapes

To explore the flexibility of Efficient Teacher, we extend it to the scenario of domain adaptive object detection

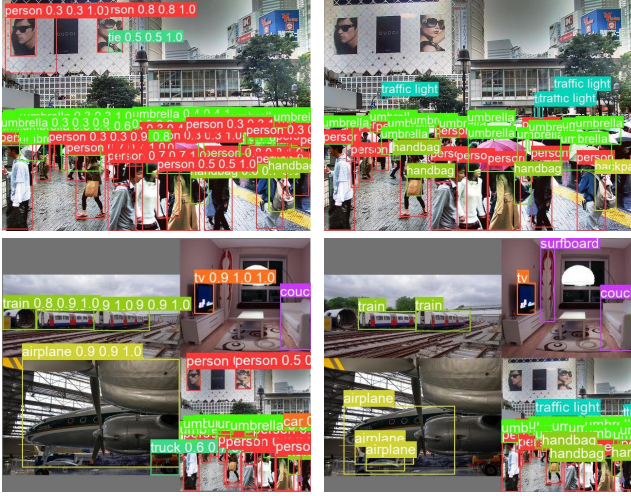


Figure 8. The upper section displays the pseudo labels and ground truth obtained using the Unbiased Teacher method. The lower section presents figures that have been augmented with Mosaic in the Efficient Teacher method. In these figures, the three floating point numbers that appear with the class name represent the pseudo label score, objectness score, and classification score, respectively.

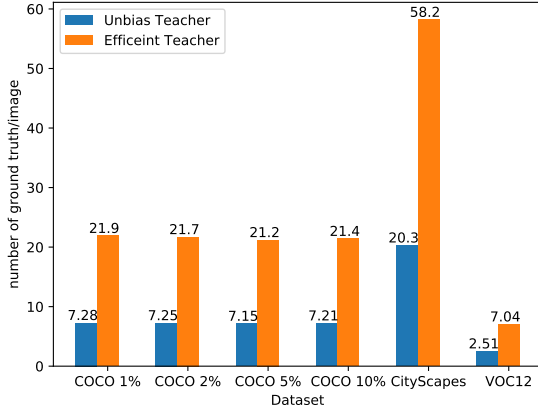


Figure 9. Number of annotation and labels after Mosaic per image among different datasets. When applying Mosaic technique, the number of labels becomes nearly three times the number of annotations, which disrupts the original label distribution.

(DAOD) [6, 17] follow the LabelMatch [4]. The experiments are mainly to demonstrate the generality of Efficient Teacher framework, which simply treats the source data as labeled data and the target data as unlabeled data.

**Dataset.** As described in Tab. 7, following the LabelMatch, We evaluate our method on these settings and compare it with the state-of-the-arts.

**Implementation Details.** The implementation is nearly the same as SSOD, and more training hyper-parameters can be

Data Split	Normal→Foggy	KITTI→CityScapes	Sim10K→CityScapes
labeled data	Cityscapes (train)	KITTI	Sim10K
unlabeled data	Cityscapes-foggy (train)	Cityscapes (train)	Cityscapes (train)
test data	Cityscapes-foggy (val)	Cityscapes (val)	Cityscapes (val)

Table 7. Three different domain shifts in DAOD, which are constructed by four different datasets, including Cityscapes [9], Cityscapes-foggy [34], KITTI [16] and Sim10k [20].

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 [44]	24.7	54.1	46.7	32.9	41.1	32.4	38.7	45.7	39.5
CVPR2020:HTCN [6]	31.6	47.9	47.5	33.2	40.9	32.3	37.1	47.4	39.8
CVPR2021:MeGA [41]	25.4	52.4	49.0	37.7	46.9	34.5	39.0	49.2	41.8
CVPR2021:UMT [10]	34.1	48.6	46.7	33.0	46.8	30.4	37.3	56.5	41.7
CVPR2022:TDD [17]	35.1	68.2	53.7	50.7	45.1	38.9	49.1	53.0	49.2
CVPR2022:LabelMatch [4]	<b>42.0</b>	62.2	55.4	45.3	<b>55.1</b>	43.5	51.5	64.1	52.4
Source only	24.2	51.9	49.0	46.5	8.5	30.3	39.2	37.2	35.9
Efficient Teacher(Ours)	39.2	<b>72.9</b>	<b>59.2</b>	<b>58.1</b>	53.0	<b>46.4</b>	<b>54.8</b>	59	<b>55.3</b>

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

Method	$AP_{50}$	network
Source only	42.2	FR+VGG
CVPR2019:SW-Faster [32]	37.9	FR+VGG
CVPR2020:GPA [44]	47.9	FR+R50
CVPR2021:MeGA [41]	43.0	FR+VGG
CVPR2022:TDD [17]	47.4	FR+VGG
CVPR2022:LabelMatch [4]	51.0	FR+VGG
ICCV2021:SimROD [29]	47.5	YOLOv5
Efficient Teacher (Ours)	<b>56.4</b>	YOLOv5

Table 9. Results of adaptation from KITTI to CityScapes. FR: Faster-RCNN.

Method	$AP_{50}$	network
Source only	36.5	FR+VGG
CVPR2019:SW-Faster [32]	40.7	FR+VGG
CVPR2020:GPA [44]	47.6	FR+R50
CVPR2021:MeGA [41]	44.8	FR+VGG
CVPR2022:TDD [17]	53.4	FR+VGG
CVPR2022:LabelMatch [4]	52.7	FR+VGG
ICCV2021:SimROD [29]	52.1	YOLOv5
Efficient Teacher(Ours)	<b>59.3</b>	YOLOv5

Table 10. Results of adaptation from S10K to CityScapes. VGG: VGG-16.

found in Tab. 11. Following previous works, we use  $AP_{50}$  as our evaluation metric, moreover, only  $AP_{50}$  on car is reported in KITTI to CityScapes and S10K to CityScapes setting because of category differences in annotations.

**Results.** We validate the DAOD results of Efficient Teacher on CityScapes, Foggy CityScapes, KITTI and S10K datasets. In the experiments, the input resolution is

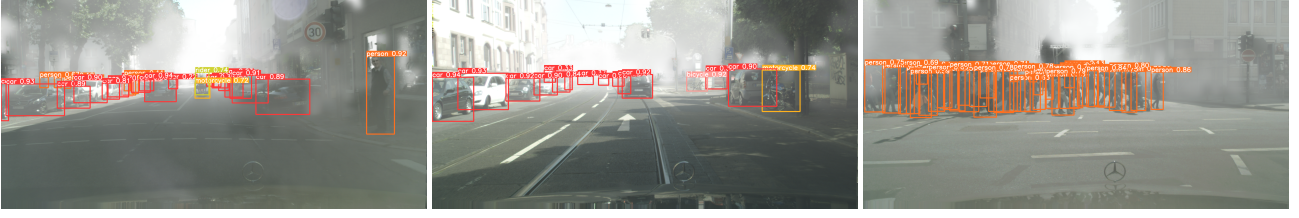


Figure 10. DAOD Performance of YOLOv5 with Efficient Teacher on Foggy Cityscapes validation dataset.

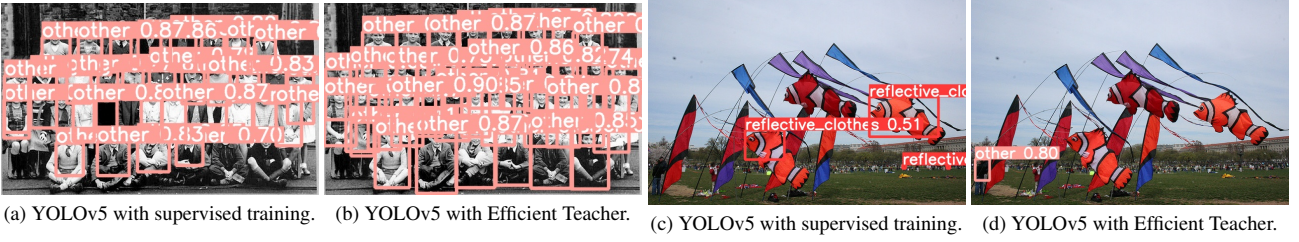


Figure 11. Comparison between supervised object detection and semi-supervised object detection on customized dataset.

960 and  $\lambda$  is 1.0. Table 8 shows Efficient Teacher outperforms all previous state-of-the-art model with a margin of 2.9  $AP_{50}$ . Table 9 and Tabel 10 shows our method ,is an effective framework for both SSOD and DAOD.

### C. Experiments on Customized Datasets

We further report the performance of Efficient Teacher on customized dataset to showcase the generality of our framework.

**Dataset.** The customized dataset contains two types of ground truth: reflective clothes and other, the relective clothes means the annotation of the reflective clothes that people wear and other means normal clothes. We only collect 3800 pictures in a few scenes while customer wish the model can be deployed in a lot of scenes, thus we take 3800 pictures as labeled data and COCO train dataset as unlabeled data for Efficient Teacher.

**Implementation Details.** We verify the generality of Efficient Teacher through inference the supervised model and semi-supervised model on COCO val dataset, results are visualized in Figure 11. Obviously, the SSOD on Efficient Teacher shows great performance improvement with fewer classification and localization error. In addition, what is interesting in the figures should be the red carp flags are misclassified as reflective clothes, which is the common phenomenon occurs by supervised training in practical applications. Efficient Teacher shows superiority of alleviating the overfitting due to small training dataset and potential to become the standard solution in application scenario.

### D. Implementation and Training Details

Our Efficient Teacher is based on YOLO-like Dense Detector, for fair comparison, we refactor the source code of YOLOv5(<https://github.com/ultralytics/yolov5>), adding utilities like unlabeled dataloader, pseudo labels online transformation and configuration system. Our code can train both supervised and semi-supervised object detection through modify a few lines of configuration. The source code will be released soon.

training setting	COCO-standard	COCO-additional	VOC	DAOD	Customized
input resolution	640	640	640	960	640
batch size for labeled data	32	32	32	16	32
batch size for unlabeled data	32	32	32	16	32
learning rate	0.01	0.01	0.01	0.01	0.01
iterations	300K	540K	72K	20K	150K
unsupervised loss weight $\lambda$	3.0	2.0	2.0	1.0	3.0
EMA rate	0.999	0.999	0.999	0.999	0.999
reliable ratio $\alpha$	0.5	0.5	0.5	0.5	0.3
pseudo label NMS score thresh	0.01	0.01	0.01	0.01	0.1
pseudo label NMS IoU thresh	0.65	0.65	0.65	0.65	0.65
multi-scale (strong augmentation)	(0.1, 1.9)	(0.1, 1.9)	(0.1, 1.9)	(0.5, 1.5)	(0.5, 1.5)
test score thresh	0.001	0.001	0.001	0.001	0.001

Table 11. 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
Mosaic	1.0	None	.
Strong Augmentation			
Process	Prob	Parameters	Descriptions
Mosaic	1.0	None	.
Horizontal Flip	0.5	None	None
Multi-Scale	1.0	ratio=(0.1, 1.9)	The short edge of image is random resized from $0.1l_{short}$ to $1.9l_{short}$ .
HSV Color Jittering	1.0	(brightness, saturation, hue) = (0.4, 0.7, 0.015)	Brightness factor is chosen uniformly form [0.6, 1.4], saturation factor is chosen uniformly from [0.3, 1.7], and hue value is chosen uniformly from [-0.015, 0.015].
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 12. Details of data augmentations.