

# Unbiased Teacher v2: Semi-supervised Object Detection for Anchor-free and Anchor-based Detectors

Yen-Cheng Liu<sup>1</sup>, Chih-Yao Ma<sup>2</sup>, Zsolt Kira<sup>1</sup>

<sup>1</sup>Georgia Institute of Technology, <sup>2</sup>Meta

{ycliu, zkira}@gatech.edu, cyma@fb.com

## Abstract

*With the recent development of Semi-Supervised Object Detection (SS-OD) techniques, object detectors can be improved by using a limited amount of labeled data and abundant unlabeled data. However, there are still two challenges that are not addressed: (1) there is no prior SS-OD work on anchor-free detectors, and (2) prior works are ineffective when pseudo-labeling bounding box regression. In this paper, we present Unbiased Teacher v2, which shows the generalization of SS-OD method to anchor-free detectors and also introduces Listen2Student mechanism for the unsupervised regression loss. Specifically, we first present a study examining the effectiveness of existing SS-OD methods on anchor-free detectors and find that they achieve much lower performance improvements under the semi-supervised setting. We also observe that box selection with centerness and the localization-based labeling used in anchor-free detectors cannot work well under the semi-supervised setting. On the other hand, our Listen2Student mechanism explicitly prevents misleading pseudo-labels in the training of bounding box regression; we specifically develop a novel pseudo-labeling selection mechanism based on the Teacher and Student’s relative uncertainties. This idea contributes to favorable improvement in the regression branch in the semi-supervised setting. Our method, which works for both anchor-free and anchor-based methods, consistently performs favorably against the state-of-the-art methods in VOC, COCO-standard, and COCO-additional.*

## 1. Introduction

Deep learning models have achieved remarkable performance on object detection tasks in recent years, though the strong performance heavily relies on training a network with abundant images with human-annotated labels. To reduce the label supervision for training object detectors, Semi-Supervised Object Detection (SS-OD) methods have been proposed to leverage only limited labeled data but more abun-

dant unlabeled data to improve performance [5, 9, 20, 26, 38]. Existing state-of-the-art SS-OD methods apply self-training techniques, which generate pseudo-labels and enforce the consistency between unlabeled data with different augmentations. Despite the significant improvement, there are still two remaining issues that are left untackled: (1) **there is no prior SS-OD work on anchor-free detectors** and (2) **prior works are ineffective in pseudo-labeling on the bounding box regression**.

First, anchor-free detectors have been recently getting more attention in the community of object detection [11, 15, 16, 29, 30, 37, 42], with the promise of achieving competitive accuracy, computational efficiency, and potential generalization to new datasets or environments [37]. In spite of these advances, existing SS-OD works [9, 20, 26] mainly focus on anchor-based detectors (e.g., Faster-RCNN [21] and SSD [19]) but do not empirically verify their effectiveness on anchor-free detectors. In fact, when we adapt recent state-of-the-art SS-OD methods to anchor-free detectors, we observe that, compared with its improvement on anchor-based models, the improvement is much smaller on anchor-free models (see Figure 1a and Table 1). With extensive analysis provided in Section 3.2, we find that *some advanced techniques performing favorably in the fully-supervised setting do not work in the semi-supervised setting with limited supervision*. For example, the centerness score becomes unreliable for box selection under the semi-supervised setting, and the localization-based labeling method is not robust to the localization noise in pseudo-labels.

Second, following the Teacher-Student framework, the existing SS-OD works [26, 38] apply an unsupervised regression loss with the pseudo-boxes generated from confidence thresholding (i.e., a threshold on the box score). However, we find that this approach inherits some potential issues that can be further addressed. For instance, (1) instead of using one single metric (e.g., box score or box IoU) to *jointly* represent the quality of four boundaries, the confidence/uncertainty of **each** boundary should be predicted individually; (2) confidence in the classification branch might not be able to reflect the quality of boundary prediction on

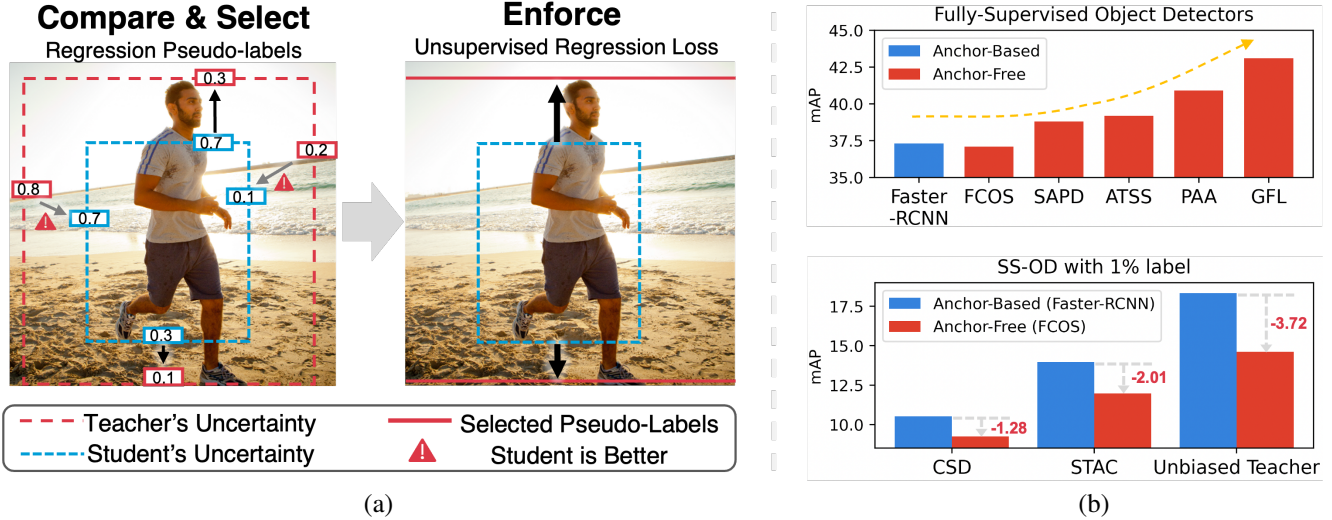


Figure 1. To improve the unsupervised regression loss, we propose (a) *Listen2Student*, which explicitly compares the prediction uncertainties between the Teacher and the Student and selects these instances where the teacher has lower uncertainty than the student. We then enforce the unsupervised regression loss on these selected regression pseudo-labels. (b) Anchor-free detectors are rapidly developed recently, while adapting the pseudo-labeling method on the anchor-free models results in less improvements compared with the anchor-based detectors.

the regression branch. Instead, we propose to predict uncertainties on the regression branch to select pseudo-labels for boundary prediction; (3) Lastly, simply relying on Teacher’s confidence/uncertainty prediction to select pseudo labels for regression cannot prevent *misleading* instances for the regression task. Instead, we propose to exploit the *relative uncertainties between the Teacher and Student* to select the boundary-level pseudo-labels, in which the Teacher has lower uncertainty than the Student. Integrating the three components, we propose *Listen2Student* to improve the unsupervised regression loss for the SS-OD tasks, as shown in Figure 1b.

We demonstrate that our proposed method achieves significant improvements compared to the state-of-the-art SS-OD methods when using *both* anchor-free and anchor-based detectors on several SS-OD benchmarks, including *COCO-standard*, *COCO-additional*, and *VOC*. We also provide ablation studies to examine the effectiveness of our *Listen2Student*. We summarize the main contributions as follows:

- We show the generalization of our proposed semi-supervised method on both anchor-based and anchor-free detectors. To the best of our knowledge, we are the first to examine the anchor-free models on SS-OD, and we identify core issues in applying SS-OD methods on anchor-free detectors.
- We explicitly remove misleading instances in regression pseudo-labels by considering *relative* uncertainty estimation from the Teacher and Student predictions. We provide analyses to verify effectiveness of our approach on anchor-free and anchor-based detectors.

- Based on our empirical study on anchor-free and anchor-based detectors, our method shows favorable improvements against the state-of-the-art methods. With the proposed method, we also bridge the performance gap between anchor-free and anchor-based detectors under the semi-supervised setting.

## 2. Related Work

**Anchor-Free Object Detectors.** The development of deep learning models has resulted in significant improvements on object detection tasks. Existing object detectors consist of anchor-based detectors [2, 17, 19, 21, 24, 32] and anchor-free detectors [11, 13, 29, 30, 40, 42]. Specifically, anchor-based detectors predict the box shift and scaling for the pre-defined anchor-boxes, and each predicted box is labeled according to its intersection-over-union (IoU) score to the ground-truth boxes. Based on label assignment (*i.e.*, assign classification labels to predicted instances) and sub-sampling of foreground-background anchor boxes, the models are then trained to perform object detection. Despite remarkable results that have been achieved, applying anchor-based detectors on new datasets requires experts to tune hyper-parameters [10] related to anchor-boxes, which limits the ability to adapt to new datasets or environments [37].

Alternatively, anchor-free models alleviate these concerns by removing the pre-defined anchor-boxes in detection models. For example, keypoint-based anchor-free detectors eliminated the need for designing a set of anchor boxes by representing a box as two corner points [13], a center point with four extreme points [40], and a center point with the box weight and height [39]. Similarly, FCOS [29] removed

the pre-defined anchor-boxes and predicted a classification score, distances to four boundaries, and a centerness score for each pixel. Several works improved the performance of the anchor-free model by proposing an adaptive sample selection [37], jointly training the centerness and classification branches with soft-labels [16], soft-selecting the pyramid levels [41], and modeling the boundary uncertainty [15]. In this paper, we use FCOS [29] as our base anchor-free model, since it is publicly available and widely used in existing anchor-free models [15, 16, 37, 41].

**Semi-Supervised Object Detection.** Semi-supervised learning (SSL) for image classification has been rapidly developed and obtained promising results in recent years. Existing SSL image classification works [1, 6, 8, 12, 22, 25, 28, 35, 36] apply input augmentations/perturbations and consistency regularization on unlabeled images to improve the model trained with the limited amount of labeled data. Inspired by these works, several semi-supervised object detection works have been proposed to exploit similar ideas to train object detectors in a semi-supervised manner. For example, CSD [9] apply a left-right consistency loss to enforce prediction consistency between horizontally flipped unlabeled images. Some other works [20, 26, 27, 38] exploit pseudo-labeling, where a model iteratively generates the pseudo-labels of unlabeled data and add the confident predictions into the training data. STAC [26] uses the limited amount of labeled data to train an object detector, which is used to generate the pseudo-labels for unlabeled data in an offline manner. To refine the quality of pseudo-labels, Instant-Teaching [38] proposes a co-rectify scheme to rectify the false prediction between two identical but independently trained models. Humble Teacher [27] applies exponential moving average (EMA) and soft pseudo-labels to improve against the model trained on labeled data only. Unbiased Teacher [20] proposes to generate the pseudo-labels in an online fashion, and the quality of pseudo-labels is further improved by addressing the pseudo-labeling bias issue. Soft-Teacher [33] proposes a simple background-weighted loss and box variance filter to improve performance against the supervised baselines. While they can improve the performance in the semi-supervised setting, existing works only present their results on anchor-based detectors. We are thus interested in investigating the generalization of the state-of-the-art methods (*i.e.*, pseudo-labeling) on anchor-free models and improving the performance of anchor-free models for semi-supervised object detection tasks.

### 3. Method

#### 3.1. Background: Semi-supervised Object Detection and Pseudo-labeling

With the goal of learning an object detector in a semi-supervised setting, we assume a set of labeled images

Table 1. **Adaption of Unbiased Teacher [20] to an anchor-free model.** The performance is degraded when applying Unbiased Teacher on the anchor-free model (FCOS).

Methods	Models	COCO-standard				100%
		0.5%	1%	5%	10%	
UT [20]	F-RCNN	14.36	18.33	26.65	29.56	37.90
UT [20]	FCOS	10.27 (-4.09)	14.61 (-3.72)	23.99 (-2.66)	28.18 (-1.38)	38.10

$D_s = \{x_i^s, y_i^s\}_{i=1}^{N_s}$  and unlabeled images  $D_u = \{x_i^u\}_{i=1}^{N_u}$  are available during training.

In order to address semi-supervised object detection, existing works [20, 26, 38] exploit the pseudo-labeling method. Specifically, this line of works contains two stages: 1) The burn-in stage and 2) the mutual learning stage. In the burn-in stage, with the available labeled data, an initial object detector is trained with the standard supervised losses,  $\mathcal{L}_{sup} = \sum_i \mathcal{L}(x_i^s, y_i^s)$ . In the mutual learning stage, the pretrained object detector is duplicated into a Student and a Teacher model initially. Then, in each training iteration, the Teacher model takes the weakly-augmented unlabeled images as input and predicts the bounding boxes, and the instances with the box score higher than a threshold  $\tau$  (*i.e.*, confidence thresholding) are selected as the pseudo-labels.

Based on the pseudo-labels and the same unlabeled image but with a stronger augmentation, the unsupervised loss  $\mathcal{L}_{unsup}$  is computed and combined with the supervised loss  $\mathcal{L}_{sup}$  to train the Student model,  $\theta_s \leftarrow \theta_s + \gamma \frac{\partial(\mathcal{L}_{sup} + \lambda_u \mathcal{L}_{unsup})}{\partial \theta_s}$ , where  $\mathcal{L}_{unsup} = \sum_i \mathcal{L}(x_i^u, \hat{y}_i^u)$ . To refine the quality of the pseudo-labels, the Teacher model weight ( $\theta_t$ ) can be further updated with the Student model weight ( $\theta_s$ ) via Exponential Moving Average (EMA) as shown in [20].

Although the existing works [20, 26, 38] based on pseudo-labeling have presented significant improvements on anchor-based detectors (*i.e.*, Faster-RCNN), it is still unclear whether such a method is applicable to anchor-free detectors. This motivates us to investigate its generalization to anchor-free detectors, and we provide our findings and show that the state-of-the-art SS-OD method is not effective as it is mostly designed for anchor-based detectors (in Section 3.2).

#### 3.2. Pseudo-labeling on Anchor-Free Detectors

We take the widely used FCOS model [29] as an example of anchor-free detector for studying SS-OD task. FCOS [29] has three major prediction branches: 1) a classifier for performing object category classification, 2) a centerness branch for indicating the probability of being the center of foreground objects, and 3) a regressor for estimating the distances to the boundaries of an object. These models usually exploit fully convolutional layers and perform pixel-wise predictions. To train the model, all pixels inside the ground-truth boxes are labeled as foreground and the remaining pixels as background, and regression loss and centerness loss are only

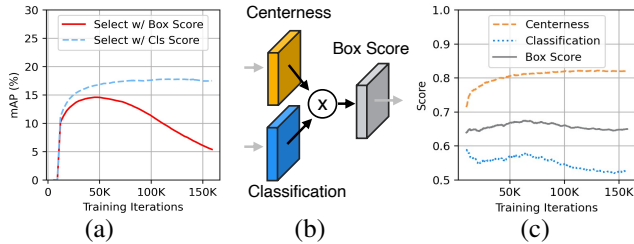


Figure 2. **Illustration of Centerness bias issue.** (a) Selecting pseudo-boxes based on box scores leads to worse results in semi-supervised learning compared with selecting based on classification scores. (b) Box scores of the anchor-free detectors [29, 37] are defined as the multiplication of the centerness scores and the classification scores, and we find that (c) box scores of pseudo-boxes is dominated by the centerness scores, which are unreliable in semi-supervised setting (see Appendix for further details).

Table 2. While box selection based on box score leads to higher detection accuracy in fully-supervised settings, it performs worse than box selection based on classification scores under semi-supervised learning settings. Fully-supervised results are from FCOS [29].

	Class. score	Box score	$\Delta$
Fully-supervised	33.50	37.10	<b>+3.60</b>
Semi-supervised	17.79	15.12	<b>-2.67</b>

enforced in these foreground instances. For more details of anchor-free detectors, please refer to the FCOS paper [29].

As shown in Figure 1b and Table 1, we observe that simply applying the existing state-of-the-art SS-OD methods [9, 20, 26] on anchor-free detectors obtains much smaller improvements compared with anchor-based detectors. We attribute this to the following two factors.

**Centerness bias issue.** As presented in Figure 2b and Table 2, we notice that selecting the pseudo-boxes based on box scores performs worse than solely relying on classification scores in the semi-supervised setting, while FCOS [29] shows using box scores leads to better results in the fully-supervised setting. We observed that this is because the box scores of some anchor-free detectors [29, 37] are defined as the multiplication of classification scores and centerness scores (see Figure 2a), and the pseudo-boxes selected based on the box scores have relatively high centerness scores but low classification scores (see Figure 2c). This reveals that the box scores are dominated by the centerness scores in the pseudo-labeling mechanism. However, with the limited amount of labels used in the training, the centerness scores are not reliable for reflecting whether a prediction is a foreground instance since there is no supervision to suppress the centerness scores for background instances in the centerness branch<sup>1</sup>. As a result, these selected high centerness pseudo-

<sup>1</sup>A similar observation was also made in Generalized Focal Loss [16].

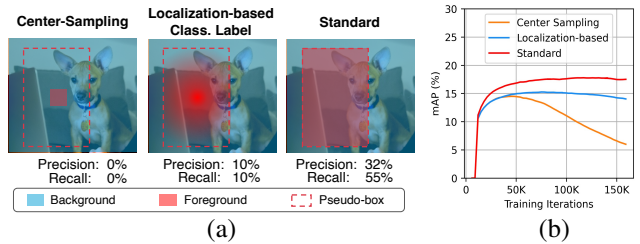


Figure 3. **Illustration of Unreliable label assignments.** (a) Existing techniques for improving fully-supervised anchor-free detectors such as Center Sampling [29] and localization-based classification labels [16] are less robust to the localization noise (e.g., box center shifted) in pseudo-boxes, and the pixel-wise recall and precision of these two techniques are lower than the standard label assignment. Thus, (b) our empirical evaluation shows that the standard label assignment leads to a better result.

Table 3. While the Center-Sampling improves the anchor-free detectors in fully-supervised setting, it degrades the performance in the semi-supervised setting. Fully-supervised results from FCOS [29].

	w/o Center-Sampling	w/ Center-Sampling	$\Delta$
Fully-supervised	37.10	38.10	<b>+1.00</b>
Semi-supervised	17.79	14.96	<b>-2.83</b>

boxes are likely to be the background instances, and adding these false-positive pseudo-boxes in the semi-supervised training degrades the effectiveness of the pseudo-labeling and also aggravates the centerness bias issue.

**Unreliable Label Assignment.** To improve the performance of the *fully-supervised* anchor-free detector, several works [16, 39] proposed to use soft classification labels, which are weighted based on the bounding box localization as presented in Figure 3a. Similarly, FCOS [29] also presented an advanced label assignment technique, *center-sampling*, which regards the instances close to the center of the object as foreground instances and improves against the model using the standard label assignment that labels all instances inside ground-truth boxes as foreground and the remaining instances as background. *Although the above techniques improve the anchor-free detectors during fully-supervised training, we found that they are not effective or even detrimental during semi-supervised training* (see Figure 3b and Table 3). We hypothesize that this is because the pseudo-boxes can have localization noise (either due to the center of the box being shifted or the box has incorrect width and height), and using *center-sampling* or the *Localization-based soft labels* makes pixel-wise predictions incorrectly labeled as either foreground (false positive) or background (false negative). For instance, as shown in Figure 3, the precision and recall of *center-sampling* is much lower than *standard* for this particular example with reasonable amount of localization noise.

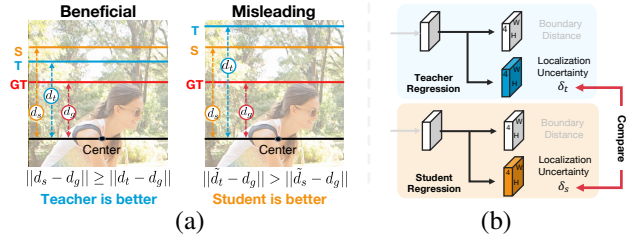


Figure 4. (a) **Beneficial/Misleading regression pseudo-labels** and (b) *Listen2Student*. (a) We categorize the regression pseudo-labels into beneficial and misleading instances, and (b) our *Listen2Student* prevents the misleading regression pseudo-labels and thus improves the regression branch.

To alleviate the centerness bias issue, we select the pseudo-boxes based on classification score only (and ignore the centerness score) in limited-supervision scenarios, since we empirically found that classification score is more reliable to represent the objectness of the predicted instances especially under limited supervision. In this way, the false-positive pseudo-labels are less likely to impede the effectiveness of pseudo-labeling and thus improve the performance of pseudo-labeling. We also train the classifier with the hard labels (*i.e.*, one-hot vector) rather than the soft labels with the box localization weighting. Finally, instead of using the *center-sampling*, we use the *standard* label assignment method, which labels all elements inside the bounding boxes as the foreground and the remaining as the background.

### 3.3. Listen2Student for Unsupervised Regression Loss

#### 3.3.1 Limitations of confidence thresholding for regression

While confidence thresholding has been demonstrated to work well in classification (image-level [25] or box-level [26, 38]), we observed solely relying on the box confidence cannot effectively remove the *misleading* instances in the box regression, and there are several reasons why it does not perform favorably: (1) First, the confidence thresholding in existing works selects pseudo-boxes based on the box scores, which only reflect the confidence of object classification in Faster-RCNN [20], and there is no explicit module estimating the confidence (or uncertainty) of regression prediction, *i.e.*, the regression branch only predicts the boundary location without any metric indicating localization uncertainty in vanilla object detectors. (2) Second, using one single score (*e.g.*, centerness or IoU score) to jointly represent the quality of four predicted boundaries is not accurate, as it is hard to obtain a pseudo-box with four **equally precise** boundaries under the limited-supervision setting. (3) Lastly, unlike pseudo-labels for discrete object categories, the real-valued regression outputs are unbounded. Selecting pseudo-boxes solely based on the Teacher’s confidence

cannot explicitly prevent *misleading* instances in the pseudo label for regression, because the Teacher can still provide a regression direction that is contradictory to the ground-truth direction. Similar observations are also found in prior works in knowledge distillation for regression tasks [3, 23].

#### 3.3.2 Listen2Student

To address the above concerns and improve the regression branch with the Teacher-Student mechanism, we aim to select the *beneficial* instances and remove *misleading* instances for the training of the regression branch. Intuitively, we develop a novel way to use *relative* prediction information between the Student and Teacher; to our knowledge this is the first instance of moving beyond using just the Teacher’s prediction information. Specifically, as shown in Figure 4, the *beneficial* instance of boundary prediction is defined as: the instance that satisfies  $\|\tilde{d}_t - d_g\| \leq \|\tilde{d}_s - d_g\|$ , where  $\tilde{d}_t$  is the Teacher’s regression prediction,  $\tilde{d}_s$  is the Student’s regression prediction, and  $d_g$  is the ground-truth regression label. As a comparison, the *misleading* instance of regression is expressed as the instance satisfying  $\|\tilde{d}_t - d_g\| > \|\tilde{d}_s - d_g\|$ .

**Uncertainty prediction for regression.** While we were hoping to use ground-truth labels  $d_g$  to decide whether the predictions from the Teacher is better or not, in reality, the ground-truth labels are not available for SS-OD. Therefore, we propose to predict the localization uncertainty, which loosely correlates with the error to the ground-truth label (*i.e.*,  $\|\tilde{d}_t - d_g\|$  and  $\|\tilde{d}_s - d_g\|$ ) for the unlabeled data. As shown in Figure 4, the localization uncertainty of each boundary prediction is derived by adding an additional branch, which has the same output size as the boundary distance regression branch. The localization uncertainty branch is jointly trained with the boundary distance branch, and we use the negative power log-likelihood loss (NPLL) [14]<sup>2</sup> as the regression loss,

$$\mathcal{L}_{reg}^{sup} = \sum_i \eta_i \left( \sum \left( \frac{(d_s - d_g)^2}{2\delta_s^2} + \frac{1}{2} \log \delta_s^2 \right) + 2 \log 2\pi \right),$$

where  $\eta_i$  is the IoU score between the predicted box and the ground-truth box, and  $\delta_s$  is the predicted uncertainty of the Student.

**Relative uncertainties for pseudo-label selection.** With the uncertainty estimation, we first loosely remove the boundaries where student has very small localization uncertainty  $\delta_s \leq \sigma_s$ . We then propose a selection mechanism which explicitly takes not only the Teacher’s localization uncertainty  $\delta_t^i$  but also the Student’s localization uncertainty  $\delta_s^i$  into account for the pseudo-label selection. By selecting the *beneficial* instances where the Teacher has lower localization uncertainty than the Student with a margin  $\sigma$ , our

<sup>2</sup> *Listen2Student* is not limited to NPLL, and other regression uncertainty estimation methods [7] are also potentially applicable.

unsupervised regression loss is thus defined as

$$\mathcal{L}_{reg}^{unsup} = \begin{cases} \sum_i ||\tilde{d}_t^i - \tilde{d}_s^i||, & \text{if } \delta_t^i + \sigma \leq \delta_s^i \\ 0, & \text{otherwise} \end{cases}, \quad (1)$$

where  $\sigma \geq 0$  is a margin between the localization uncertainties of Teacher and Student. Note that the unsupervised regression loss is computed in the boundary level rather than the box level, so some boundaries of a box are used to computed unsupervised regression loss while the others are not.

The core idea of this mechanism is that the Teacher should only guide the Student with the instances that the Teacher has lower uncertainty than the student, as it indicates that the Teacher has a potentially lower error. By contrast, for the instances, which the Teacher has higher uncertainty than the Student, we should not enforce the loss, as the Teacher is likely to predict worse than the Student and thus mislead the Student for these instances. Based on this selection mechanism, we can explicitly prevent gradients from *misleading* instances from degrading the performance of the regression branch. Our regression branch can thus gradually be refined and obtain more accurate boundary prediction. It is worth noting that the localization uncertainty branch is an individual branch and only used in the training stage, thus introducing no additional computation during inference.

## 4. Experiments

### 4.1. Settings and Implementation Details

**Experimental Settings.** We follow the experimental settings presented in the existing semi-supervised object detection works [20, 26]. Specifically, we use MS-COCO [18] and PASCAL VOC [4] and examine our proposed method on three experimental scenarios, *COCO-standard*, *COCO-additional*, and *VOC*. For *COCO-standard*, we randomly sample 0.5, 1, 2, 5, and 10% labeled training data as our labeled set, and the remaining data as the unlabeled set. For *COCO-additional*, we use *COCO2017-labeled* as labeled set and *COCO2017-unlabeled* as the unlabeled set. We evaluate on *COCO2017-val* for both *COCO-standard* and *COCO-additional* as in previous works. As for *VOC*, *VOC2007-trainval* is used as the labeled set, and *VOC2012-trainval* and *COCO20cls* are used as the unlabeled set. All trained models in *VOC* experiment are evaluated on *VOC2007-test*.

**Model Architecture.** In order to examine the effectiveness of anchor-free models for semi-supervised object detection, we chose FCOS [29] as our base anchor-free models since it is widely adopted in existing anchor-free works [15, 16, 37, 41]. As the existing works mainly focus on anchor-based models and use Faster-RCNN [9, 20, 26] or SSD [9], we also adapt the existing SS-OD methods [9, 20, 26] to the anchor-free model (*e.g.*, FCOS).

**Implementation Details.** Our implementation is based on Detectron2 [31]. To train our model, we use SGD optimizer with the learning rate 0.01, and each batch contains 8 labeled images and 8 unlabeled images unless specified. We use the unsupervised loss weight  $\lambda_u = 3.0$  and classification threshold  $\tau = 0.5$ , and we set  $\sigma = 0.1$  as the margin between localization uncertainties of Teacher and Student and  $\sigma_s = 0.5$ . We adapt the data augmentation used in Unbiased Teacher and applied the scale jittering used in SoftTeacher [33] without using any geometric augmentation during training, as we empirically find that the scale jittering leads to a significant improvement. More details are listed in the supplementary material.

### 4.2. Results on Anchor-free Detector

**COCO-standard.** We adapt three anchor-based methods, CSD [9], STAC [26], and Unbiased Teacher [20], to the anchor-free models, and each method was ran five times and their means and variances are reported, as presented in Table 4. Our model consistently performs favorably against the baseline methods under different degrees of supervision, and the improvement gap is larger when the level of supervision is lower. Our experiments on VOC and COCO-additional also result in a similar trend as well (see Appendix for experimental results).

### 4.3. Results on Anchor-based Detector

In addition to the results on the anchor-free model, we are also interested whether our proposed method can generalize to different types of object detectors. Specifically, we apply our unsupervised regression loss on Unbiased Teacher and modify the regression branch to predict the localization uncertainty with an additional branch as we did in Section 3.3. We examine our *Listen2Student* on the Faster-RCNN for *COCO-standard*, *VOC*, and *COCO-additional* as follows.

**COCO-standard.** As presented in Table 5, compared with the state-of-the-art SS-OD methods [20, 27, 33], our method obtains higher mAP under the cases where 0.5% to 10% data are labeled. Under different batch sizes, we could maintain the improvement gap against existing SS-OD methods and further improve the performance to 35.08 mAP under *COCO-standard* 10% case. In addition, we also find that the performance gap between the anchor-free and anchor-based detectors is reduced by using our framework, and this verifies the generalization of our proposed *Listen2Student* to both anchor-free and anchor-based detectors.

**VOC and COCO-additional.** To verify whether our framework can improve the object detector trained with the unlabeled set, we also consider *VOC* in Table 7 and *COCO-additional* in Table 8. With *VOC07* used as the labeled set, our model can leverage *VOC12* to achieve 56.87 mAP, and using *VOC12+COCO20cls* as the unlabeled set can further improve the model and achieve 58.08mAP. On the

Table 4. Experimental results of the anchor-free model (FCOS-ResNet50) on *COCO-standard*. \* We reimplement and adapt to FCOS-ResNet50. We randomly sample labeled data and run each method 5 times, and we report the mean and standard deviation for each result. We used 8 labeled images and 8 unlabeled images for all results presented in this table.

	Anchor-free detectors on COCO-standard				
	0.5%	1%	2%	5%	10%
Supervised	5.42 ± 0.01	8.43 ± 0.03	11.97 ± 0.03	17.01 ± 0.01	20.98 ± 0.01
CSD [9]*	5.76 ± 0.55 (+0.34)	9.23 ± 0.08 (+0.80)	12.53 ± 0.04 (+0.56)	18.09 ± 0.08 (+1.08)	22.06 ± 0.01 (+1.08)
STAC [26]*	8.79 ± 0.12 (+3.37)	11.97 ± 0.12 (+3.54)	15.50 ± 0.16 (+3.53)	20.36 ± 0.05 (+3.35)	24.31 ± 0.02 (+3.33)
Unbiased Teacher [20]*	10.27 ± 0.13 (+4.85)	14.61 ± 0.10 (+6.18)	18.70 ± 0.21 (+6.73)	23.99 ± 0.12 (+6.98)	28.18 ± 0.01 (+7.20)
<b>Ours</b>	<b>16.25 ± 0.18 (+10.83)</b>	<b>22.71 ± 0.42 (+14.28)</b>	<b>26.03 ± 0.12 (+14.06)</b>	<b>30.08 ± 0.04 (+13.07)</b>	<b>32.61 ± 0.03 (+11.63)</b>

Table 5. Experimental results of anchor-based models (FasterRCNN-ResNet50) on *COCO-standard*. For a fair comparison, we make the batch size consistent to the baseline methods. †: using labeled/unlabeled batch size 32/32, \*: using batch size labeled/unlabeled batch size 8/40, and rest of the results using batch size 8/8. We randomly sample labeled data and run each method 5 times, and we report the mean and standard deviation for each result.

	Anchor-based detectors on COCO-standard				
	0.5%	1%	2%	5%	10%
Supervised	6.83 ± 0.15	9.05 ± 0.16	12.70 ± 0.15	18.47 ± 0.22	23.86 ± 0.81
CSD [9]	7.41 ± 0.21 (+0.58)	10.51 ± 0.06 (+1.46)	13.93 ± 0.12 (+1.23)	18.63 ± 0.07 (+0.16)	22.46 ± 0.08 (-1.40)
STAC [26]	9.78 ± 0.53 (+2.95)	13.97 ± 0.35 (+4.92)	18.25 ± 0.25 (+5.55)	24.38 ± 0.12 (+5.86)	28.64 ± 0.21 (+4.78)
Humble Teacher [27]	-	16.96 ± 0.38 (+7.91)	21.72 ± 0.24 (+9.02)	27.70 ± 0.15 (+9.23)	31.61 ± 0.28 (+7.74)
Instant Teaching [38]	-	18.05 ± 0.15 (+9.00)	22.45 ± 0.15 (+9.75)	26.75 ± 0.05 (+8.28)	30.40 ± 0.05 (+6.54)
Unbiased Teacher [20]	14.36 ± 0.09 (+7.53)	18.33 ± 0.19 (+9.28)	22.23 ± 0.21 (+9.53)	26.65 ± 0.31 (+8.18)	29.56 ± 0.24 (+5.70)
ISMT [34]	-	18.88 ± 0.74 (+9.83)	22.43 ± 0.56 (+9.73)	26.37 ± 0.24 (+7.90)	30.53 ± 0.52 (+6.67)
<b>Ours</b>	<b>17.51 ± 0.24 (+10.68)</b>	<b>21.84 ± 0.13 (+12.79)</b>	<b>26.14 ± 0.01 (+13.44)</b>	<b>30.06 ± 0.14 (+11.59)</b>	<b>33.50 ± 0.03 (+9.64)</b>
SoftTeacher [33]*	-	20.46 ± 0.39 (+11.41)	-	30.74 ± 0.08 (+12.27)	34.04 ± 0.14 (+10.18)
<b>Ours*</b>	<b>21.02 ± 0.49 (+14.19)</b>	<b>24.79 ± 0.30 (+15.74)</b>	<b>28.23 ± 0.05 (+15.53)</b>	<b>32.05 ± 0.04 (+13.58)</b>	<b>35.02 ± 0.02 (+11.16)</b>
Unbiased Teacher [20]†	16.94 ± 0.23 (+10.11)	20.75 ± 0.12 (+11.72)	24.30 ± 0.07 (+11.60)	28.27 ± 0.11 (+9.80)	31.50 ± 0.10 (+7.64)
<b>Ours†</b>	<b>21.26 ± 0.21 (+14.43)</b>	<b>25.40 ± 0.36 (+16.35)</b>	<b>28.37 ± 0.03 (+15.67)</b>	<b>31.85 ± 0.09 (+13.38)</b>	<b>35.08 ± 0.02 (+11.22)</b>

other hand, with the *COCO2017-unlabeled* set, our model can perform favorably against the object detector trained on *COCO2017-train* and achieve 44.75 mAP. Note that we train our model for 720k iterations and do not tune the inference threshold (same as SoftTeacher). Training the model longer or tuning the inference threshold can potentially further improve the performance. These results confirm the effectiveness of our framework on improving the existing object detector using the extra unlabeled images.

#### 4.4. Effectiveness of Unsupervised Regression Loss

We compare the methods including 1) our proposed *Listen2Student* 2) No unsupervised regression loss, and 3) using confidence thresholding and enforcing L1 loss, as used in existing works [26, 38]. To further understand how these methods contribute to the improvement of bounding box regression, we provide an mAP breakdown from AP55 to

AP95 of each method in Table 6. It is worth noting that we only change the unsupervised regression loss across these methods and keep the remaining objective functions and modifications the same across all variants.

We observe that, although the confidence thresholding can improve the easier evaluation metrics (e.g., AP55), it cannot improve or even degrades the results on stricter evaluation metrics (e.g., AP95). This shows that simply using the confidence thresholding cannot prevent misleading pseudo-labels from degrading the performance on extremely precise boundary predictions. In contrast, our *Listen2Student* shows consistent improvements on all evaluation metrics and leads to favorable results, especially on these stricter evaluation metrics. This empirically confirms that our *Listen2Student* contributes to the more precise bounding box prediction, as our *Listen2Student* enforces the boundary-wise unsupervised regression loss, which exploits the pseudo-labels derived by

Table 6. Average precision (AP) breakdown of unsupervised regression methods. We also report the absolute improvement of each unsupervised regression loss method against the model without the unsupervised regression loss.

	AP55	AP60	AP65	AP 70	AP75	AP80	AP85	AP90	AP95
No regression	29.71	27.34	24.64	21.38	17.55	13.27	8.33	3.45	0.35
Confidence Thresholding	30.60 <b>+0.89</b>	28.19 <b>+0.85</b>	25.07 <b>+0.43</b>	21.93 <b>+0.55</b>	17.96 <b>+0.41</b>	13.32 <b>+0.05</b>	8.22 <b>-0.11</b>	3.12 <b>-0.33</b>	0.32 <b>-0.03</b>
<i>Listen2Student</i> (Ours)	30.78 <b>+1.07</b>	28.59 <b>+1.25</b>	26.19 <b>+1.56</b>	23.05 <b>+1.67</b>	19.64 <b>+2.09</b>	15.61 <b>+2.34</b>	10.47 <b>+2.14</b>	5.06 <b>+1.61</b>	0.58 <b>+0.23</b>

Table 7. Results of the **Anchor-based model (Faster-RCNN)** on VOC.

Methods	Labeled	Unlabeled	$AP^{50}$	$AP^{50:95}$
Supervised			76.70	43.60
STAC [26]			77.45	44.64
ISMT [34]			77.23	46.23
Instant-Teaching [38]	VOC07	VOC12	79.20	50.00
Humble Teacher [27]			80.94	53.04
Unbiased Teacher [20]			80.51	54.48
Ours			<b>81.29</b>	<b>56.87</b>
STAC [26]			79.08	46.01
ISMT [34]			77.75	49.59
Instant-Teaching [38]	VOC07	VOC12	79.00	50.80.
Humble Teacher [27]		+	81.29	54.41
Unbiased Teacher [20]		<i>COCO20cls</i>	81.71	55.79
Ours			<b>82.04</b>	<b>58.08</b>

comparing the uncertainty estimation of each boundary prediction.

**Limitations and future works.** Although we have shown the improvement and generalization on anchor-free and anchor-based detectors, applying SSOD methods on a large-scale unlabeled dataset (*e.g.*, OpenImage) remains a challenge. We also find that the localization uncertainty estimation for boundary prediction leaves room for improvement to be integrated with the relative thresholding mechanism. There are other challenges such as unseen objects in the unlabeled dataset or domain shift between datasets. While these topics are not our focus in this paper, they are worth exploring in future research.

## 5. Conclusion

In this paper, we examined the existing SS-OD methods on anchor-free models and presented the SS-OD benchmarks on anchor-free detectors. By identifying and addressing the core issues that existed in the pseudo-labeling method on anchor-free detectors, our method can improve against the state-of-the-art methods. We further pre-

Table 8. Results of the **Anchor-based model (Faster-RCNN)** on *COCO-additional*. \*We adapt the scale jitter used in Soft-Teacher [33] to Unbiased Teacher and it leads to significant improvement.

Methods	$mAP$
Supervised	40.90
CSD [9]	38.52
STAC [26]	39.21
Humble Teacher [27]	42.37
Unbiased Teacher* [20]	44.06
SoftTeacher [33]	44.50
Ours	<b>44.75</b>

sented *Listen2Student*, a novel method that uses *relative* Teacher/Student uncertainties to explicitly prevent the misleading regression pseudo-labels and select beneficial regression pseudo-labels in a boundary-wise manner. This enables the regression branch to benefit from the use of unlabeled images. In the experiment sections, we examine each method in three different SS-OD tasks and present consistent improvements. We also provide an extensive study to verify the effectiveness and generalization of our proposed *Listen2Student* mechanism on both anchor-free and anchor-based detectors.

Concerning negative societal impacts, we think it is essential to be aware that there exists the risk that object detection techniques (not just our method) are used in surveillance systems. Also, as this line of works relies on low-labeled data for the model training, this aggravates the risk of data bias toward historically disadvantaged groups.

## 6. Acknowledgments

Yen-Cheng Liu and Zsolt Kira were partly supported by DARPA’s Learning with Less Labels (LwLL) program under agreement HR0011-18-S-0044, as part of their affiliation with Georgia Tech.



## References

- [1] 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 (NeurIPS)*, pages 5049–5059, 2019. [3](#)
- [2] Zhaowei Cai, Quanfu Fan, Rogerio S Feris, and Nuno Vasconcelos. A unified multi-scale deep convolutional neural network for fast object detection. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016. [2](#)
- [3] Guobin Chen, Wongun Choi, Xiang Yu, Tony Han, and Manmohan Chandraker. Learning efficient object detection models with knowledge distillation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. [5](#)
- [4] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision (IJCV)*, 88(2):303–338, 2010. [6](#)
- [5] Jiyang Gao, Jiang Wang, Shengyang Dai, Li-Jia Li, and Ram Nevatia. Note-rcnn: Noise tolerant ensemble rcnn for semi-supervised object detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019. [1](#)
- [6] Hongyu Guo, Yongyi Mao, and Richong Zhang. Mixup as locally linear out-of-manifold regularization. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 33, pages 3714–3722, 2019. [3](#)
- [7] Yihui He, Chenchen Zhu, Jianren Wang, Marios Savvides, and Xiangyu Zhang. Bounding box regression with uncertainty for accurate object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2888–2897, 2019. [5](#)
- [8] Dan Hendrycks, Norman Mu, Ekin D. Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. AugMix: A simple data processing method to improve robustness and uncertainty. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020. [3](#)
- [9] Jisoo Jeong, Seungeui Lee, Jeesoo Kim, and Nojun Kwak. Consistency-based semi-supervised learning for object detection. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. [1](#), [3](#), [4](#), [6](#), [7](#), [8](#)
- [10] Licheng Jiao, Fan Zhang, Fang Liu, Shuyuan Yang, Lingling Li, Zhixi Feng, and Rong Qu. A survey of deep learning-based object detection. *IEEE Access*, 7:128837–128868, 2019. [2](#)
- [11] 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. [1](#), [2](#)
- [12] Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017. [3](#)
- [13] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018. [2](#)
- [14] Youngwan Lee, Joong-won Hwang, Hyung-Il Kim, Kimin Yun, and Joungyoul Park. Localization uncertainty estimation for anchor-free object detection. *arXiv preprint arXiv:2006.15607*, 2020. [5](#)
- [15] Xiang Li, Wenhai Wang, Xiaolin Hu, Jun Li, Jinhui Tang, and Jian Yang. Generalized focal loss v2: Learning reliable localization quality estimation for dense object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. [1](#), [3](#), [6](#)
- [16] Xiang Li, Wenhai Wang, Lijun Wu, Shuo Chen, Xiaolin Hu, Jun Li, Jinhui Tang, and Jian Yang. Generalized focal loss: Learning qualified and distributed bounding boxes for dense object detection. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. [1](#), [3](#), [4](#), [6](#)
- [17] 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 (CVPR)*, pages 2980–2988, 2017. [2](#)
- [18] 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 *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014. [6](#)
- [19] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision (ECCV)*, pages 21–37. Springer, 2016. [1](#), [2](#)
- [20] 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 *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021. [1](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#)
- [21] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems (NeurIPS)*, pages 91–99, 2015. [1](#), [2](#)
- [22] Mehdi Sajjadi, Mehran Javanmardi, and Tolga Tasdizen. Regularization with stochastic transformations and perturbations for deep semi-supervised learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 1163–1171, 2016. [3](#)
- [23] Muhamad Risqi U Saputra, Pedro PB de Gusmao, Yasin Almalioglu, Andrew Markham, and Niki Trigoni. Distilling knowledge from a deep pose regressor network. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019. [5](#)
- [24] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors with online hard example mining. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. [2](#)
- [25] 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 (NeurIPS)*, 2020. [3](#), [5](#)
- [26] 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, 3, 4, 5, 6, 7, 8
- [27] 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. 3, 6, 7, 8
- [28] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *Advances in neural information processing systems (NeurIPS)*, pages 1195–1204, 2017. 3
- [29] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019. 1, 2, 3, 4, 6
- [30] Jiaqi Wang, Kai Chen, Shuo Yang, Chen Change Loy, and Dahua Lin. Region proposal by guided anchoring. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1, 2
- [31] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2. <https://github.com/facebookresearch/detectron2>, 2019. 6
- [32] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [33] 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. 3, 6, 7, 8
- [34] 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/CVF Conference on Computer Vision and Pattern Recognition*, pages 5941–5950, 2021. 7, 8
- [35] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 6023–6032, 2019. 3
- [36] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *Proc. International Conference on Learning Representations (ICLR)*, 2018. 3
- [37] 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 Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 1, 2, 3, 4, 6
- [38] 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 Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 3, 5, 7, 8
- [39] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. In *arXiv preprint arXiv:1904.07850*, 2019. 2, 4
- [40] Xingyi Zhou, Jiacheng Zhuo, and Philipp Krahenbuhl. Bottom-up object detection by grouping extreme and center points. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [41] Chenchen Zhu, Fangyi Chen, Zhiqiang Shen, and Marios Savvides. Soft anchor-point object detection. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 3, 6
- [42] Chenchen Zhu, Yihui He, and Marios Savvides. Feature selective anchor-free module for single-shot object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1, 2