

1 Thanks for your constructive reviews. We tried our best to properly respond to all the inquiries from the reviewers.

2 1. Common Response

3 **Comparison with SSM:** For a practical usage of a model, the stopping criteria must be defined properly. SSM
4 continues the training until no more unlabeled data are included in the training set (Algorithm 1 in [1]). In Fig.4 of [1],
5 the performance improves initially as unlabeled samples are added but it starts to degrade as more samples are added.
6 However, the reported score of the SSM is not from an objective stopping criterion but the peak performance during
7 the entire iterations. With this setting, the influx of the data from out-of-distribution and incorrectly labeled samples
8 cannot be prevented. As we all know that the performance of SSM should not be measured with this setting, we are not
9 sure that the performance of a detector combined with SSM would get better. To check this, we tried to implement
10 SSM in SSD. However, many details are missing in SSM and the learning parameters of single-stage detector and
11 two-stage detector are different. Its intense time-cost and huge hyper-parameter space makes it difficult to implement
12 SSM properly. In our work, we just wanted to present a representative self-learning method.

13 **Dataset:** When the unlabeled samples originate from out-of-class distribution, the performance of any semi-supervised
14 learning method usually degrades. This is a well known limitation of SSL [2] and we wanted to show the performance
15 of our method in the in-class distribution setting as other works.

16 Also, we acknowledge that using a significantly huge and unrefined data would be an ideal setting. However, this
17 paper has its focus on applying the consistency-based semi-supervised approach to the object detection task and the
18 experimental settings have been inspired from those in the paper of SSM. We would like to mention on the dataset-related
19 research direction that some reviewers pointed out in an additional paragraph in the final version.

20 2. Response to each Reviewer

21 **Reviewer 1:** In consistency regularization, various candidates of loss function can be used. We simply chose the Π
22 model, which uses L_2 loss, as a baseline and we will explain this more in detail in the final version.

23 Particularly in the case of using con-l in SSD300, the performance is quite good even without using BE. We think that
24 it is because the number of samples with flat areas is relatively small for a small input resolution and this helps the
25 regression process while con-c disturbs the effect of con-l.

26 We have experimented with SSD300 model for the effect of not using CR on labeled samples. SSD300 (CR on unlabeled
27 samples only) scored 72.1mAP in semi-supervised learning using VOC07(labeled) and VOC12(unlabeled) and it is
28 slightly worse than the base CSD (72.3 mAP).

29 We will search for more related works including the ones the reviewer mentioned and will try to explore a link between
30 the works and Fig. 1 in our paper.

31 **Reviewer 2:** Table 2 shows the results of experiments with 20 COCO classes and the entire 80 classes. As the reviewer
32 mentioned, we would like to modify our claim (line 246-248) more clearly as follows: The reason why using 20 class
33 categories as unlabeled data outperforms the other case is because the ratio of *labeled/unlabeled class mismatch* is
34 much smaller. This is similar to the experimental results of Fig. 2 in [2].

35 **Reviewer 3:** As mentioned in the paper, the self-training method is an iterative method. So it is time-consuming and
36 computationally intensive. In addition, the mAP score is reported wrongly in SSM, as mentioned above. Meanwhile,
37 our CSD is a new type of semi-supervised learning method for object detection that can be trained with an additional
38 loss reflecting consistency between pairs of images.

39 **Reviewer 4:** Both hard negative mining and focal loss require a label in their training. To apply these to the consistency
40 loss of unlabeled data, they should be modified appropriately. Therefore, we propose BE, which is applicable in the
41 semi-supervised learning environment. It seems that the reviewer considered the weighted loss using the background
42 score, which we think can be effective that can be another new contribution. We will apply it in our future research.

43 It is important to clarify the relationship between boxes by image perturbation to calculate the consistency loss. As the
44 reviewer mentioned, the performance is expected to be improved by diversity if some limited perturbation schemes are
45 applied. Although some of the perturbation methods the reviewer mentioned have minor problems (e.g image wrapping
46 can cut off an object, which is not representable by a single bounding box), we agree that additional perturbation
47 methods will probably improve the performance. However, flipping is the most simple method which always guarantees
48 the one-to-one correspondence of given boxes. In the discussion, we will complement the possibility of using other
49 perturbation methods.

50 References

51 [1] Keze Wang, Xiaopeng Yan, Dongyu Zhang, Lei Zhang, and Liang Lin. Towards human-machine cooperation: Self-supervised
52 sample mining for object detection. In *CVPR*, pages 1605–1613, 2018.

53 [2] Avital Oliver, Augustus Odena, Colin A Raffel, Ekin Dogus Cubuk, and Ian Goodfellow. Realistic evaluation of deep semi-
54 supervised learning algorithms. In *NIPS*, pages 3235–3246, 2018.