

A Simple Semi-Supervised Learning Framework for Object Detection

Kihyuk Sohn*, Zizhao Zhang*, Chun-Liang Li,
Han Zhang, Chen-Yu Lee, and Tomas Pfister

Google Cloud AI, Google Brain

{ksohn, zizhao, chunliang, zhanghan, chenyllee, tpfister}@google.com

Abstract. Semi-supervised learning (SSL) has promising potential for improving the predictive performance of machine learning models using unlabeled data. There has been remarkable progress, but the scope of demonstration in SSL has been limited to image classification tasks. In this paper, we propose STAC, a simple yet effective SSL framework for visual object detection along with a data augmentation strategy. STAC deploys highly confident pseudo labels of localized objects from an unlabeled image and updates the model by enforcing consistency via strong augmentations. We propose new experimental protocols to evaluate performance of semi-supervised object detection using MS-COCO and demonstrate the efficacy of STAC on both MS-COCO and VOC07. On VOC07, STAC improves the $AP^{0.5}$ from 76.30 to 79.08; on MS-COCO, STAC demonstrates $2\times$ higher data efficiency by achieving 24.38 mAP using only 5% labeled data than supervised baseline that marks 23.86% using 10% labeled data. The code is available at https://github.com/google-research/ssl_detection/.

Keywords: Semi-supervised learning, object detection, self-training, augmentation consistency

1 Introduction

Semi-supervised learning (SSL) has received growing attention in recent years as it provides means of using unlabeled data to improve model performance when large-scale annotated data is not available. A popular class of SSL methods is based on Consistency-based Self-Training [27,38,25,44,54,35,4,58,3,59,49]. The key idea is to first generate the artificial labels for the unlabeled data and train the model to predict these artificial labels when feeding the unlabeled data with semanticity-preserving stochastic augmentations. The artificial label can either be a one-hot prediction (hard) or the model’s predictive distribution (soft). The other pillar for the success of SSL is from advancements in data augmentations. Data augmentations improve the robustness of deep neural networks [48,24] and has been shown to be particularly effective for consistency-based self-training [58,3,59,49]. The augmentation strategy spans from a manual

* Equal contribution.

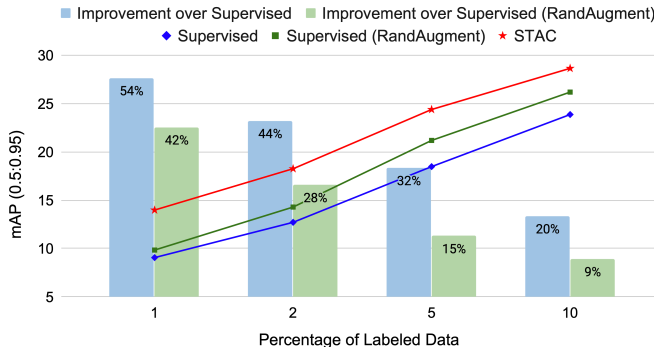


Fig. 1: The proposed semi-supervised learning framework for object detection, STAC, consistently improves upon supervised baselines and those with data augmentation using different amount of labeled training data on MS-COCO [30].

combination of basic image transformations, such as rotation, translation, flipping, or color jittering, to neural image synthesis [68,21,62] and policies learned by reinforcement learning [6,69]. Lately, complex data augmentation strategies, such as RandAugment [7] or CTAugment [3], have turned out to be powerful for SSL on image classification [58,3,49,59].

While having made remarkable progress, SSL methods have been mostly applied to image classification, whose labeling cost is relatively cheaper compared to other important problems in computer vision, such as object detection. Due to its expensive labeling cost, object detection demands a higher level of label efficiency, necessitating the development of strong SSL methods. On the other hand, the majority of existing works on object detection has focused on training a stronger [47,28,29] and faster [14,41,8] detector given sufficient amount of annotated data, such as MS-COCO [30]. Few existing works on SSL for object detection [53,34,43] rely on additional context, such as categorical similarities of objects or temporal consistency from video.

In this work, we leverage lessons learned from deep SSL on image classification to tackle SSL for object detection. To this end, we propose a SSL framework for object detection that combines self-training (via pseudo label) [46,33] and consistency regularization based on the strong data augmentations [6,7,69]. Inspired by the framework in Noisy-Student [59], our system contains two stages of training. In the first stage, we train an object detector (e.g., Faster RCNN [41]) using all labeled data until convergence. The trained detector is then used to predict bounding boxes and class labels of localized objects for unlabeled images as shown in Figure 2. Then, we apply confidence-based box filtering to each predicted box (after non-maximum suppression) with high threshold value to obtain pseudo labels with high precision, inspired by the design of FixMatch [49]. In the second stage, the strong data augmentations are applied to each unlabeled image and the model is trained with labeled data and unlabeled data with its corresponding pseudo labels generated in the first stage. Encouraged by RandAugment [7] and its successful adaptation to SSL [58,49] and object

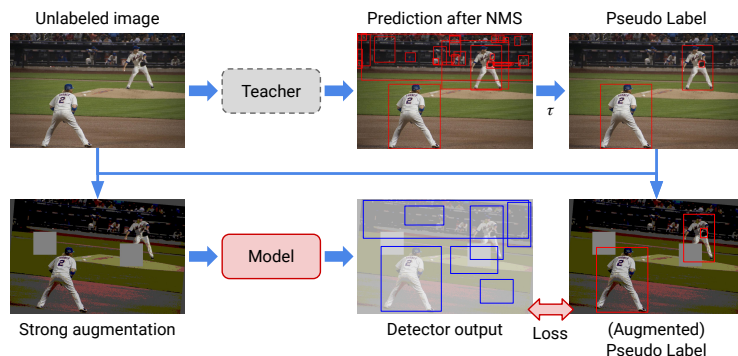


Fig. 2: The proposed SSL framework for object detection. We generate pseudo labels (i.e., bounding boxes and their class labels) for unlabeled data using test-time inference, including non-maximum suppression (NMS) [15], of the teacher model trained with labeled data. We compute unsupervised loss with respect to pseudo labels whose confidence scores are above a threshold τ . The strong augmentations are applied for augmentation consistency during the model training. Target boxes are augmented when global geometric transformations are used.

detection [69], we design our augmentation strategy for object detection, which consists of global color transformation, global or box-level [69] geometric transformations, and Cutout [10].

We test the efficacy of STAC on public datasets: MS-COCO [30] and PASCAL VOC [13]. We design new experimental protocols using MS-COCO dataset to evaluate the semi-supervised performance of object detection. We use 1, 2, 5 and 10% of labeled data as labeled sets and the remainder as unlabeled sets to evaluate the effectiveness of SSL methods in the low-label regime. In addition, following [37,52], we evaluate using all labeled data as the labeled set and additional unlabeled data provided by MS-COCO as the unlabeled set. Following [23], we use trainval of VOC07 as the labeled set and that of VOC12 with or without unlabeled data of MS-COCO as unlabeled sets. While being simple, STAC brings significant gain in mAPs: 18.47 to 24.38 on 5% protocol, 23.86 to 28.64 on 10% protocol as in Figure 1, and 42.60 to 46.01 on PASCAL VOC.

The contribution of this paper is as follows:

1. We develop STAC, a SSL framework for object detection that seamlessly extends the class of state-of-the-art SSL methods for classification based on self-training and augmentation-driven consistency regularization.
2. STAC is simple and introduces only two new hyperparameters: the confidence threshold τ and the unsupervised loss weight λ_u , which do not require an extensive additional effort for tuning.
3. We propose new experimental protocols for SSL object detection using MS-COCO and demonstrate the efficacy of STAC on MS-COCO and PASCAL VOC in Faster RCNN framework.

2 Related Work

Object detection is a fundamental computer vision task and has been extensively studied in the literature [15,14,41,17,28,5,39,40,31,29]. Popular object detection frameworks include Region-based CNN (RCNN) [15,14,41,17,28], YOLO [39,40], SSD [31], etc [26,51,11]. The progress made by existing works is mainly on training a stronger or faster object detector given sufficient amount of annotated data. There is growing interest in improving detectors using unlabeled training data through a semi-supervised object detection framework [53,34]. Before deep learning, the idea has been explored by [42]. Recently, [23] proposes a consistency-based semi-supervised object detection method, which enforces the consistent prediction of an unlabeled image and its flipped counterpart. Their method requires a more sophisticated Jensen-Shannon Divergence for consistency regularization computation. Similar ideas to consistency regularization have also been studied in the active learning settings for object detection [55]. [52] introduces a self-supervised proposal learning module to learn context-aware and noise-robust proposal features from unlabeled data. [37] proposes data distillation that generates labels by ensembling predictions of multiple transformations of unlabeled data. We argue that stronger semi-supervised detectors require further investigation of unsupervised objectives and data augmentations.

Semi-supervised learning (SSL) for image classification has been dramatically improved recently. Consistency regularization becomes one of the popular approaches among recent methods [2,45,25,63,58] and inspires [23] on object detection. The idea is to enforce the model to generate consistent predictions across label-preserving data augmentations. Some exemplars include Mean-Teacher [54], UDA [58], and MixMatch [4]. Another popular class of SSL is pseudo labeling [27,2], which can be viewed as a hard version of consistency regularization: the model is performing self-training to generate pseudo labels of unlabeled data and thereby train randomly-augmented unlabeled data to match the respective pseudo labels (i.e. being consistent in predictions of the same unlabeled example). How to use pseudo labels is critical to the success of SSL. For instance, Noisy-Student [59] demonstrates an iterative teacher-student framework that repeats the process of labeling assignments using a teacher model and then training a larger student model. This method achieves state-of-the-art performance on ImageNet classification by leveraging extra unlabeled images in the wild. FixMatch [49] demonstrates a simple algorithm which outperforms previous approaches and establishes state-of-the-art performance, especially on diverse small labeled data regimes. The key idea behind FixMatch is matching the prediction of the strongly-augmented unlabeled data to the pseudo label of the weakly-augmented counterpart when the model confidence on the weakly-augmented one is high. In light of the success of these methods, this paper exploits the effective usage of pseudo labeling and pseudo boxes as well as data augmentations to improve object detectors.

Data augmentations are critical to improve model generalization and robustness [6,7,19,69,64,67,10,12,20], especially gradually become a major impetus on semi-supervised learning [4,3,58,49]. Finding appropriate color transformations

and geometric transformations of input spaces has been shown to be critical to improve generalization [6,19]. However, most augmentations are mainly studied in image classification. The complexity of data augmentations for object detection is much higher than image classification [69], since global geometric transformations of data affect bounding box annotations. Some works have presented augmentation techniques for supervised object detection, such as MixUp [64,65], CutMix [61], or augmentation strategy learning [69]. The recent consistency-based SSL object detection method [23] utilizes global horizontal flipping (weak augmentation) to construct the consistency loss. To the best of our knowledge, the impact of intensive data augmentations on semi-supervised object detection has not been well studied.

3 Methodology

3.1 Background: Unsupervised Loss in SSL

Formulating an unsupervised loss that leverages unlabeled data is the key in SSL. Many advancements in SSL for classification rely on some forms of *consistency regularization* [27,38,25,44,54,35,4,58,3,59,49]. Inspired by a comparison in [49], we provide a unified view of consistency regularization for image classification. For K -way classification, the consistency regularization is written as follows:

$$\ell_u = \sum_{x \in \mathcal{X}} w(x) \ell(q(x), p(x; \theta)) \quad (1)$$

where $x \in \mathcal{X}$ is an image, $p, q: \mathcal{X} \rightarrow [0, 1]^K$ map x into a $(K-1)$ -simplex, and $w: \mathcal{X} \rightarrow \{0, 1\}$ maps x into a binary value. $\ell(\cdot, \cdot)$ measures a distance between two vectors. Typical choices include L_2 distance and cross entropy. Here, p represents the prediction of the model parameterized by θ , q is the prediction target, and w is the weight that determines the contribution of x to the loss. As an example, pseudo labeling [27] has the following configurations:

$$q(x) = \text{ONE_HOT}(\arg \max(p(x; \theta))), w(x) = 1 \text{ if } \max(p(x; \theta)) \geq \tau \quad (2)$$

We refer readers to the supplementary material for configurations of a comprehensive list of SSL methods. State-of-the-art SSL algorithms, such as Unsupervised Data Augmentation (UDA) [58] or FixMatch [49], apply strong data augmentation \mathcal{A} , such as RandAugment [58] or CTAugment [3], to the model prediction $p(\mathcal{A}(x); \theta)$ for improved robustness. Noisy-Student [59] applies diverse forms of stochastic noise to the model prediction, including input augmentations via RandAugment, and network augmentations via dropout [50] and stochastic depth [22]. While sharing similarities on the model prediction, they differ in q that generates the prediction target as detailed in Appendix C. Besides the use of soft or hard targets, Different from (2) and many aforementioned algorithms, Noisy-Student employs a “teacher” network other than $p(\cdot, \theta)$ to generate pseudo labels $q(x)$. Note that the teacher network is independent of the model at training and this gives a scalability and flexibility on the choice of network architectures or optimization such as learning schedules.

3.2 STAC: SSL for Object Detection

We develop a novel SSL framework for object detection, called **STAC**, based on the **Self-Training** (via pseudo label) and the **Augmentation driven Consistency regularization**. First, we adopt a stage-wise training of Noisy-Student [59] for its scalability and flexibility. This involves at least two stages of training, where in the first stage, we train a teacher model using all available labeled data, and in the second stage, we train STAC using both labeled and unlabeled data. Second, we use a threshold with a high value for the confidence-based thresholding inspired by FixMatch [49] to control the quality of pseudo labels in object detection, which is comprised of bounding boxes and their class labels. The steps for training STAC are summarized as follows:

1. **Train a teacher model** on available labeled images.
2. **Generate pseudo labels** of unlabeled images (i.e., bounding boxes and their class labels) using the trained teacher model.
3. **Apply strong data augmentations** to unlabeled images, and augment pseudo labels (i.e. bounding boxes) correspondingly when global geometric transformations are applied.
4. **Compute unsupervised loss** and supervised loss to train a detector.

Training a Teacher Model. We develop our formulation based on the Faster RCNN [41] as it has been one of the most representative detection framework. Faster RCNN has a classifier (CLS) and a region proposal network (RPN) heads on top of the shared backbone network. Each head has two modules, namely region classifiers (e.g., a $(K+1)$ -way classifier for the CLS head or a binary classifier for the RPN head) and bounding box regressors. We present the supervised and unsupervised losses of the Faster RCNN for the RPN head for simplicity. The supervised loss is written as follows:

$$\begin{aligned} \ell_s(x, \mathbf{p}^*, \mathbf{t}^*) &= \sum_b \ell_{s,b}(x, \mathbf{p}_b^*, t_b^*) \\ &= \sum_b \left[\frac{1}{N_{\text{cls}}} \sum_i \mathcal{L}_{\text{cls}}(p_i, p_{i,b}^*) + \frac{\lambda}{N_{\text{reg}}} \sum_i p_{i,b}^* \mathcal{L}_{\text{reg}}(t_i, t_b^*) \right] \end{aligned} \quad (3)$$

where b is an index of the ground-truth bounding box and i is an index of an anchor. p_i is the predictive probability of an anchor being positive, t_i is the 4-dimensional coordinates of an anchor. $p_{i,b}^*$ is the binary label of an anchor with respect to the box b , t_b^* is the ground-truth box coordinates of the box b . To train an RPN, $p_{i,b}^*$ needs to be determined for all anchors, box pairs. Note that we define a loss per box for presentation clarity, which is slightly different from that in [41].

Generating Pseudo Labels. We perform a test-time inference of the object detector from the teacher model to generate pseudo labels. That being said, the pseudo label generation involves not only the forward pass of the backbone,

RPN and CLS networks, but also the post-processing such as non-maximum suppression (NMS). This is different from conventional approaches for classification where the confidence score is computed from the raw predictive probability. We use the score of each returned bounding box after NMS, which aggregates the prediction probabilities of anchor boxes. Using box predictions after NMS has an advantage over using raw predictions (before NMS) since it removes repetitive detection. However, this does not filter out boxes at wrong locations as visualized in Figure 2 and Figure 5a.

Unsupervised Loss. When given an unlabeled image x and set of predicted bounding boxes and their confidence scores, we determine $q_{i,b}^*$, a binary label of an anchor i with respect to the pseudo box b , for all anchor, box pairs. Let s_b^* be the box coordinates of pseudo box b . The unsupervised loss of STAC is written as follows:

$$\begin{aligned} \ell_u(x, \mathbf{q}^*, \mathbf{s}^*) &= \sum_b w_b(x) \ell_{u,b}(x, \mathbf{q}_b^*, \mathbf{s}_b^*) \\ &= \sum_b w_b(x) \left[\frac{1}{N_{\text{cls}}} \sum_i \mathcal{L}_{\text{cls}}(p_i, q_{i,b}^*) + \frac{\lambda}{N_{\text{reg}}} \sum_i q_{i,b}^* \mathcal{L}_{\text{reg}}(t_i, s_b^*) \right] \end{aligned} \quad (4)$$

where $w_b(x) = 1$ if the confidence score of the predicted box b is higher than the threshold value τ and 0 otherwise. Decomposing the loss formulation of the Faster RCNN into a sum of losses for individual boxes makes the conversion from classification (Equation (1)) to detection (Equation (4)) much more transparent. Also note that the unsupervised loss is masked *per box* instead of per image, which is well aligned with our intuition. Overall, the RPN is trained by jointly minimizing two losses as follows:

$$\ell = \ell_s(x_s, \mathbf{p}^*, \mathbf{t}^*) + \lambda_u \ell_u(\mathcal{A}(x_u, \mathbf{s}^*), \mathbf{q}^*) \quad (5)$$

where \mathcal{A} is a strong data augmentation applied to an unlabeled image. Since some transformation operations are not invariant to the box coordinates (e.g., global geometric transformation [69]), the augmentation operator \mathcal{A} is applied on the pseudo box coordinates \mathbf{s}^* as well.

The loss formulation of STAC introduces two hyperparameters τ and λ_u . In the experiments, we find that $\tau = 0.9$ and $\lambda_u \in [1, 2]$ work well. Note that the consistency-based SSL object detection method in [23] requires sophisticated three-staged weighting schedule of λ_u that includes temporal ramp-up and ramp-down. On the contrary, our system demonstrates effective performance with a simple constant weighting schedule because our framework enforces the consistency using a strong data augmentation strategy.

Data Augmentation Strategy The key factor for the success of consistency-based SSL methods, such as UDA [58] or FixMatch [49], is a strong data augmentation. While the augmentation strategy for supervised and semi-supervised image classification has been extensively studied [6, 7, 3, 58, 49], not much effort has been made yet for object detection. We extend the RandAugment for object detection used in [6] using the augmentation search space recently proposed

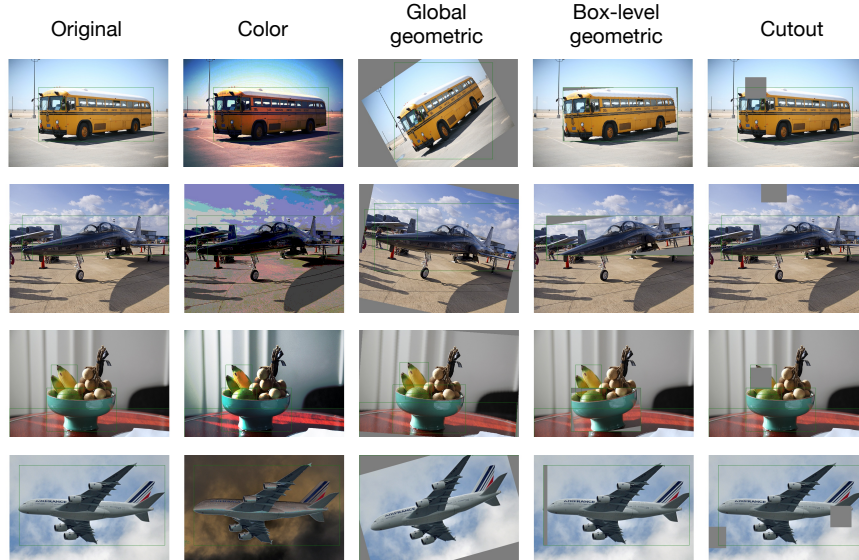


Fig. 3: Visualization of different types of augmentation strategies. From left to right: original image, color transformation, global geometric transformation, box-level geometric transformation, box-level geometric transformation, and Cutout.

by [69] (e.g., box-level transformation) along with the Cutout [10]. For completeness, we describe the list of transformation operations below. Each operation has a magnitude that decides the augmentation degree of strength.¹

1. Global color transformation (C): Color transformation operations in [7] and the suggested ranges of magnitude for each op are used.
2. Global geometric transformation (G): Geometric transformation operations in [7], namely, x - y translation, rotation, and x - y shear, are used.²
3. Box-level transformation [69] (B): Three transformation operations from global geometric transformations are used, but with smaller magnitude ranges.³

For each image, we apply transformation operations in sequence as follows. First, we apply one of the operations sampled from C. Second, we apply one of the operations sampled from either G or B. Finally, we apply Cutout at multiple random locations⁴ of a whole image to prevent a trivial solution when applied exclusively inside the bounding box. We visualize transformed images with aforementioned augmentation strategies in Figure 3.

¹ The range of degrees is empirically chosen without tuning.

² The translation range in percentage is $[-10\%, 10\%]$ of image widths or heights. The rotation and shear ranges are $[-30\%, 30\%]$ in degrees.

³ The translation range in percentage is $[-5\%, 5\%]$ of image widths or heights. The rotation and shear range is $[-10\%, 10\%]$ in degree.

⁴ The number of Cutout regions is sampled from $[1, 5]$, and the region size is sampled from $[0\%, 20\%]$ of the short edge of the applied image.

4 Experiments

We test the efficacy of our proposed method on MS-COCO [30], which is one of the most popular public benchmarks for object detection. MS-COCO contains more than 118k labeled images and 850k labeled object instances from 80 object categories for training. In addition, there are 123k unlabeled images that can be used for semi-supervised learning. We experiment two SSL settings. First, we randomly sample 1, 2, 5 and 10% of labeled training data as a labeled set and use the rest of labeled training data as an unlabeled set. For these experiments, we create 5 data folds. 1% protocol contains approximately 1.2k labeled images randomly selected from the labeled set of MS-COCO. 2% protocol contains additional ~ 1.2 k images and 5, 10% protocol datasets are constructed in a similar way. Second, following [52], we use an entire labeled training data as a labeled set and additional unlabeled data as an unlabeled set. Note that the first protocol tests the efficacy of STAC when only few labeled examples are available, while the second protocol evaluates the potential to improve the state-of-the-art object detector with unlabeled data in addition to already a large-scale labeled data. We report the mAP over 80 classes.

We also test on PASCAL VOC [13] following [23]. The trainval set of VOC07, containing 5,011 images from 20 object categories, is used as a labeled training data, and 11,540 images from the trainval set of VOC12 are used for an unlabeled training data. The detection performance is evaluated on the test set of VOC07 and mAP at IoU of 0.5 ($AP^{0.5}$) is reported in addition to the MS-COCO metric.

4.1 Implementation Details

Our implementation is based on the Faster RCNN and FPN library of Tensorpack [56]. We use ResNet-50 [18] backbone for our object detector models. Unless otherwise stated, the network weights are initialized by the ImageNet-pretrained model⁵ at all stages of training.

Since the training of the object detector is quite involved, we stay with the default learning settings for all our experiments other than the learning schedule. Most of our experiments are conducted using the `quick` learning schedule⁶ with an exception for 100% MS-COCO protocol.⁷ We find that the model’s performance is benefited significantly by longer training when more labeled training data and more complex data augmentation strategies are used. STAC introduces two new hyperparameters τ for the confidence threshold and λ_u for the unsupervised loss. We use $\tau = 0.9$ and $\lambda_u = 2$ for all experiments except for the 100% protocol of MS-COCO, where we lower threshold $\tau = 0.5$ to increase the recall of pseudo labels. We refer readers to Appendix A for complete learning settings.

⁵ <http://models.tensorpack.com/FasterRCNN/ImageNet-R50-AlignPadding.npz>

⁶ Please refer to Section 5.1 for the definition of different learning schedules.

⁷ <https://github.com/tensorpack/tensorpack/tree/master/examples/FasterRCNN#results>

⁸ We note that the number from [23] using ResNet101 with R-FCN while all the results from our implementation use ResNet50 with FPN.

Methods	1% COCO	2% COCO	5% COCO	10% COCO	100% COCO
Supervised	9.05±0.16	12.70±0.15	18.47±0.22	23.86±0.81	37.63
Supervised [†]	9.83±0.23	14.28±0.22	21.18±0.20	26.18±0.12	39.48
STAC	13.97±0.35	18.25±0.25	24.38±0.12	28.64±0.21	39.21

Table 1: Comparison in mAPs for different methods on MS-COCO. We report the mean and standard deviation over 5 data folds for 1, 2, 5 and 10% protocols. “Supervised” refers to models trained on labeled data only, which then are used to provide pseudo labels for STAC. We train STAC with the C+{B,G}+Cutout augmentation for unlabeled data. Models with [†] are trained with the same augmentation strategy, but only with labeled data. See Section 4.2 for more details.

Methods	VOC07	VOC07(AP ^{0.5})
Supervised	42.60	76.30
Supervised [†]	43.40	78.21
STAC (+VOC12)	44.64	77.45
STAC (+VOC12 & COCO)	46.01	79.08
[23] (+VOC12 & COCO)	-	75.1 ⁸

Table 2: Comparison in mAPs for different methods on VOC07. We report both mAPs at IoU=0.5:0.95, a standard metric for MS-COCO, as well as at IoU=0.5 (AP^{0.5}), since AP^{0.5} is a saturated metric as pointed out by [5]. For STAC, we follow [23] to have different level of unlabeled sources, including VOC12 and the subset of MS-COCO data with the same classes as PASCAL VOC.

4.2 Results

Since deep semi-supervised learning of visual object detectors has not been widely studied yet, we mainly compare STAC with the supervised models (i.e., models trained with labeled data only) for various experimental protocols using different data augmentation strategies. Table 1 summarizes the results. For 1, 2, 5 and 10% protocols, we train models with a quick learning schedule and report mAPs averaged over 5 data folds and their standard deviation. For 100% protocol, we employ standard with 3× longer learning schedule and report a single mAP value for each model.

Firstly, we confirm the findings of [7] with varying amount of labeled training data that the RandAugment improves the supervised learning performance of a detector by a significant margin, 2.71 mAP at 5% protocol, 2.32 mAP at 10% protocol, and 1.85 mAP for 100% protocol, upon the supervised baselines with default data augmentation of resizing and horizontal flipping.

STAC further improves the performance upon stronger supervised models. We find it to be particularly effective for protocols with small labeled training data, showing 5.91 mAP improvement at 5% protocol and 4.78 mAP at 10% protocol. Interestingly, STAC is proven to be at least 2× more data efficient than the baseline models for both 5% (24.36 for STAC v.s. 23.86 for supervised model with 10% labeled training data) and 10% protocols (28.56 for STAC v.s.

Augmentation	-	C	C+{G,B}	C+{G,B}+Cutout
5% MS-COCO (quick)	18.67	20.13	20.78	21.16
10% MS-COCO (quick)	24.05	25.26	25.92	26.34
10% MS-COCO (standard)	19.74	21.40	24.24	24.65
100% MS-COCO (standard)	37.42	37.22	36.39	36.12
100% MS-COCO (standard, 2×)	37.88	38.91	38.73	38.57
100% MS-COCO (standard, 3×)	37.63	39.33	39.75	39.48

Table 3: mAPs of supervised models trained with different augmentation and learning schedules. We test on a single fold of 5% and 10% protocols. See Section 5.1 for more details. Bold text indicates the best number in each row.

28.63 for the supervised model with 20% labeled training data). For the 100% protocol, STAC achieves 39.21 mAP. This improves upon the baseline (37.63 mAP), but falls short of the supervised model with a strong data augmentation (39.48 mAP). We hypothesize that the pseudo label training benefits by a larger amount of unlabeled data relative to the size of labeled data and study its effectiveness with respect to the scale of unlabeled data in Section 5.

We have a similar finding for experiments on PASCAL VOC. In Table 2, the mAP of the supervised models increases from 42.6 to 43.4, and $AP^{0.5}$ increases from 76.30 to 78.21. A large-scale unlabeled data from VOC12 and MS-COCO further improves the performance, achieving 46.01 mAP and 79.08 $AP^{0.5}$.

5 Ablation Study

We perform ablation study on the key components of STAC. The study analyzes the impact on the detector performance of 1) different data augmentations and learning schedule strategies, 2) different sizes of unlabeled sets, 3) the hyperparameters λ_u , coefficient for unsupervised loss, and τ , confidence threshold, and 4) quality of pseudo labels and their impact on the proposed STAC.

5.1 Data Augmentation and Learning Schedule

In this section, we evaluate the performance of supervised detector models with different data augmentation strategies and learning rate schedules while varying the amount of training data. We consider different combinations of augmentation modules, including the default augmentations of horizontal image flip, color only (C), color followed by geometric or box-level transforms (C+{G,B}), and the one followed by Cutout (C+{G,B}+Cutout). For {G,B}, we sample randomly and uniformly between geometric and box-level transform modules for each image. We consider different learning schedules, including **quick**, **standard**, and **standard** $[n]\times$ (standard setting with $[n]$ times longer training). While the number of weight updates are the same, the **quick** schedule uses lower resolution image as an input and smaller batch size for training.

The summary results are provided in Table 3. With small amount of labeled training data, we observe an increasing positive impact on detector performance

with more complex (thus stronger) augmentation strategies. The trend holds true with the `standard` schedule, but we find that the `quick` schedule is beneficial on the low-labeled data regime due to its fast training and less issue of overfitting. On the other hand, we observe that the network significantly underfits with our augmentation strategies when all labeled data is used for training. For example, with 100% labeled data, we achieve even lower mAP of 36.12 with `C+{G,B}+Cutout` strategy than that of 37.42 with default augmentations. We find that the issue can be alleviated by longer training. Moreover, while the performance with default augmentations saturates and starts to decrease as it is trained longer, the models with strong data augmentation start to outperform, demonstrating their effectiveness on training with large-scale labeled data.

STAC contains two key components: self-training and strong data augmentation. We also verify the importance of data augmentation in , which is in line with recent findings in SSL for image classification [49]. We evaluate the performance of STAC with the default augmentations (horizontal flip). On a single fold of 10% protocol, we observe a good improvement in mAP upon baseline model (from 24.05 to 26.27), but the gain is not as significant as STAC (29.00). On 100% protocol, we observe slight decrease in performance when trained with self-training only (from 37.63 to 37.57), while STAC achieves 39.21 in mAP.

5.2 Size of Unlabeled Data

While the importance of large-scale labeled data for supervised learning has been broadly studied and emphasized [9,57,30], the importance on the scale of unlabeled data for semi-supervised learning has been often overlooked [36]. In this study, we highlight the importance of large-scale unlabeled data in the context of semi-supervised object detector learning. We experiment with 5% and 10% labeled data of MS-COCO while varying the amount of unlabeled data from 1, 2, 4, and 8 times to that of the labeled data.

The summary results are given in Table 4. While there still exists the improvement in mAPs when STAC is trained with small amount of unlabeled data, the gain is less significant compared to that of supervised model with strong data augmentation. We observe clearly from Table 4 that STAC benefits from the larger amount of unlabeled training data. We make a similar observation from experiments on PASCAL VOC in Table 2, where the $AP^{0.5}$ of STAC trained using trainval of VOC12 as unlabeled data achieves 77.45, which is lower than that of supervised model with strong augmentations (78.21). On the other hand, STAC trained with large amount of unlabeled data by combining VOC12 and MS-COCO achieves 79.08 $AP^{0.5}$. This analysis may explain the slightly lower mAP of STAC for 100% protocol of MS-COCO than that of the supervised model with strong data augmentation since the size of available unlabeled data is roughly the same as that of the labeled data.

5.3 Hyperparameters λ_u and τ

We study the impact of λ_u , a regularization coefficient for unsupervised loss, and τ , the confidence threshold. Specifically, we test the STAC with different

Unlab. Size	Sup.	Sup. [†]	1×	2×	4×	8×	Full
5% MS-COCO	18.67	21.16	19.81	20.79	22.09	23.14	24.49
10% MS-COCO	24.05	26.34	25.38	26.52	27.33	27.95	29.00

Table 4: mAPs of STAC trained with varying amount of unlabeled data. $[n] \times$ refers that the amount of unlabeled data is $[n]$ times larger than that of labeled data. We test on a single fold of 5% and 10% protocols.

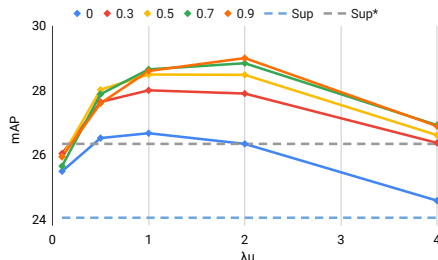


Fig. 4: mAPs of STAC with different values of $\lambda_u \in \{0.1, 0.5, 1, 2, 4\}$ and $\tau \in \{0, 0.3, 0.5, 0.7, 0.9\}$. We test on a single fold of 10% protocol. Different colors represent mAPs of models with different τ values. “Sup” represents the mAP of supervised model with default augmentations and “Sup*” represents that with $C+\{G,B\}+Cutout$.

values of $\lambda_u \in \{0.1, 0.5, 1, 2, 4\}$ and $\tau \in \{0, 0.3, 0.5, 0.7, 0.9\}$ on a single fold of 10% protocol. The summary results are provided in Figure 4. Firstly, the best performance of STAC is obtained when $\lambda_u = 2$ and $\tau = 0.9$. We observe that the performance of STAC deteriorates when λ_u is too large (> 2) or too small (< 0.5), but it improves upon strong baseline consistently for $\lambda_u \in [1, 2]$. When there is no confidence-based box filtering, the gain of STAC, if any, is marginal over the strong baseline. This is because lots of predicted boxes are indeed inaccurate, as shown in Figure 5a. Using larger value of τ allows to have pseudo box labels with higher precision (i.e., remaining boxes whose confidence is higher than τ are accurate), as in Figure 5e. However, if τ becomes too large, one would get a lower recall (e.g., bounding box at sofa in Figure 5c is filtered out in Figure 5d). Figure 4 shows that the high precision (i.e., larger value of τ) is preferred to high recall (i.e., smaller value of τ) on 10% protocol.

5.4 Quality of Pseudo Labels

One intriguing question is whether the semi-supervised performance of the model improves with pseudo labels of higher quality. To validate the hypothesis, we train two additional STAC models for 10% protocol, where models are provided pseudo labels predicted by two different supervised models trained with 5% and 100% labeled data, whose mAPs are 18.67 and 37.63, respectively. Note that the STAC on 10% protocol achieves 29.00 mAP. STAC trained with less accurate pseudo labels achieves only 24.25 mAP, while the one with more accurate pseudo labels achieves 30.30 mAP, confirming the importance of pseudo label quality.

Inspired by this observation, we increase the augmentation strength to train the teacher model in order to get better pseudo labels, expecting a further im-

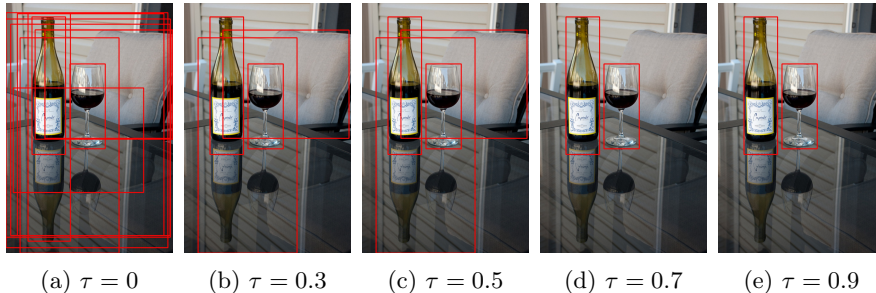


Fig. 5: Visualization of predicted bounding boxes whose confidences are larger than τ for unlabeled data. Larger value of τ results in higher precision (e.g., remaining boxes after thresholding detect objects accurately), but lower recall (e.g., detected box at sofa is removed when $\tau \geq 0.7$).

Protocol	Augmentation	-	C	C+{G,B}	C+{G,B}+Cutout
5% MS-COCO	supervised	18.67	20.13	20.78	21.16
	STAC	24.49	25.01	24.70	25.12
10% MS-COCO	supervised	24.05	25.26	25.92	26.34
	STAC	29.00	28.97	28.41	28.81

Table 5: mAPs of supervised models and STAC tested on a single fold of 5% and 10% protocols. We first train supervised models with different augmentation strategies (first row of each protocol), and pseudo labels generated from each supervised model are used to train STAC models (second row of each protocol) accordingly.

provement for STAC. To this end, we train STAC using different sets of pseudo labels that are provided by the supervised models trained with different data augmentation schemes. As in Table 5, the performance of supervised models vary from mAP of 18.67 to 21.16 with 5% labeled data and from 24.05 to 26.34 with 10% labeled data. We observe an improvement in mAP by using more accurate pseudo labels on 5% protocol, but the gain is not as substantial. We also do not observe a clear correlation between the accuracy of pseudo label and the performance of STAC on 10% protocol. While STAC brings a significant gain in mAP using pseudo labels, our results suggest that the incremental improvement on the quality of pseudo labels may not bring in a significant extra benefit.

6 Discussion and Conclusion

While SSL for classification has made significant strides, not much effort has been put to date for detection that demands a higher level of label-efficient training. We propose a simple (introducing only two hyperparameters that are easy to tune) and effective ($2\times$ label efficiency in low-label regime) SSL framework for object detection by leveraging lessons from SSL methods for classification. The

simplicity of our method will provide a flexibility for further development towards solving SSL for object detection.

The proposed framework is amenable to many variations, including using soft labels for classification loss, other detector frameworks than Faster RCNN, and other data augmentation strategies. While STAC demonstrates an impressive performance gain already without taking confirmation bias [66,1] issue into account, it could be problematic when using a detection framework with a stronger form of hard negative mining [47,29] because noisy pseudo labels can be over-used. Further investigation in learning with noisy labels, confidence calibration, and uncertainty estimation in the context of object detection are few important topics to further enhance the performance of SSL object detection.

Acknowledgment

We thank Qizhe Xie, Ekin D. Cubuk, Sercan Arik, Minh-Thang Luong, David Berthelot, Tsung-Yi Lin, Quoc V. Le, Samuel Schulter for their comments.

References

1. Arazo, E., Ortego, D., Albert, P., O'Connor, N.E., McGuinness, K.: Pseudo-labeling and confirmation bias in deep semi-supervised learning. arXiv preprint arXiv:1908.02983 (2019) 15
2. Bachman, P., Alsharif, O., Precup, D.: Learning with pseudo-ensembles. In: NeurIPS (2014) 4
3. Berthelot, D., Carlini, N., Cubuk, E.D., Kurakin, A., Sohn, K., Zhang, H., Raffel, C.: Remixmatch: Semi-supervised learning with distribution matching and augmentation anchoring. In: ICLR (2020) 1, 2, 4, 5, 7, 22, 23
4. Berthelot, D., Carlini, N., Goodfellow, I., Papernot, N., Oliver, A., Raffel, C.A.: Mixmatch: A holistic approach to semi-supervised learning. In: NeurIPS (2019) 1, 4, 5, 23
5. Cai, Z., Vasconcelos, N.: Cascade r-cnn: Delving into high quality object detection. In: CVPR (2018) 4, 10
6. Cubuk, E.D., Zoph, B., Mane, D., Vasudevan, V., Le, Q.V.: Autoaugment: Learning augmentation strategies from data. In: CVPR (2019) 2, 4, 5, 7
7. Cubuk, E.D., Zoph, B., Shlens, J., Le, Q.V.: Randaugment: Practical data augmentation with no separate search. arXiv preprint arXiv:1909.13719 (2019) 2, 4, 7, 8, 10, 22
8. Dai, J., Li, Y., He, K., Sun, J.: R-fcn: Object detection via region-based fully convolutional networks. In: NIPS (2016) 2
9. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: CVPR (2009) 12
10. DeVries, T., Taylor, G.W.: Improved regularization of convolutional neural networks with cutout. arXiv preprint arXiv:1708.04552 (2017) 3, 4, 8, 22
11. Du, X., Lin, T.Y., Jin, P., Ghiasi, G., Tan, M., Cui, Y., Le, Q.V., Song, X.: Spinenet: Learning scale-permuted backbone for recognition and localization. arXiv preprint arXiv:1912.05027 (2019) 4
12. Dwibedi, D., Misra, I., Hebert, M.: Cut, paste and learn: Surprisingly easy synthesis for instance detection. In: ICCV (2017) 4

13. Everingham, M., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. *IJCV* **88**(2), 303–338 (2010) [3](#), [9](#)
14. Girshick, R.: Fast r-cnn. In: *ICCV* (2015) [2](#), [4](#)
15. Girshick, R., Donahue, J., Darrell, T., Malik, J.: Rich feature hierarchies for accurate object detection and semantic segmentation. In: *CVPR* (2014) [3](#), [4](#)
16. Grandvalet, Y., Bengio, Y.: Semi-supervised learning by entropy minimization. In: *Advances in neural information processing systems*. pp. 529–536 (2005) [21](#)
17. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: *ICCV* (2017) [4](#)
18. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *CVPR* (2016) [9](#)
19. Hendrycks, D., Mu, N., Cubuk, E.D., Zoph, B., Gilmer, J., Lakshminarayanan, B.: Augmix: A simple data processing method to improve robustness and uncertainty. *arXiv preprint arXiv:1912.02781* (2019) [4](#), [5](#)
20. Ho, D., Liang, E., Stoica, I., Abbeel, P., Chen, X.: Population based augmentation: Efficient learning of augmentation policy schedules. *arXiv preprint arXiv:1905.05393* (2019) [4](#)
21. Hoffman, J., Tzeng, E., Park, T., Zhu, J.Y., Isola, P., Saenko, K., Efros, A.A., Darrell, T.: Cycada: Cycle-consistent adversarial domain adaptation. *arXiv preprint arXiv:1711.03213* (2017) [2](#)
22. Huang, G., Sun, Y., Liu, Z., Sedra, D., Weinberger, K.Q.: Deep networks with stochastic depth. In: *ECCV* (2016) [5](#)
23. Jeong, J., Lee, S., Kim, J., Kwak, N.: Consistency-based semi-supervised learning for object detection. In: *NeurIPS* (2019) [3](#), [4](#), [5](#), [7](#), [9](#), [10](#)
24. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: *NeurIPS* (2012) [1](#)
25. Laine, S., Aila, T.: Temporal ensembling for semi-supervised learning. In: *ICLR* (2017) [1](#), [4](#), [5](#), [21](#)
26. Law, H., Deng, J.: Cornernet: Detecting objects as paired keypoints. In: *ECCV* (2018) [4](#)
27. Lee, D.H.: Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In: *ICML Workshops* (2013) [1](#), [4](#), [5](#), [21](#)
28. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: *CVPR* (2017) [2](#), [4](#)
29. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: *ICCV* (2017) [2](#), [4](#), [15](#)
30. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: *ECCV* (2014) [2](#), [3](#), [9](#), [12](#)
31. Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C.: Ssd: Single shot multibox detector. In: *ECCV* (2016) [4](#)
32. McClosky, D., Charniak, E., Johnson, M.: Effective self-training for parsing. In: *Proceedings of the main conference on human language technology conference of the North American Chapter of the Association of Computational Linguistics*. pp. 152–159. Association for Computational Linguistics (2006) [20](#)
33. McLachlan, G.J.: Iterative reclassification procedure for constructing an asymptotically optimal rule of allocation in discriminant analysis. *Journal of the American Statistical Association* **70**(350), 365–369 (1975) [2](#)
34. Misra, I., Shrivastava, A., Hebert, M.: Watch and learn: Semi-supervised learning for object detectors from video. In: *CVPR* (2015) [2](#), [4](#)

35. Miyato, T., Maeda, S.i., Ishii, S., Koyama, M.: Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *T-PAMI* (2018) [1](#), [5](#), [21](#)
36. Oliver, A., Odena, A., Raffel, C.A., Cubuk, E.D., Goodfellow, I.: Realistic evaluation of deep semi-supervised learning algorithms. In: *NeurIPS* (2018) [12](#)
37. Radosavovic, I., Dollár, P., Girshick, R., Gkioxari, G., He, K.: Data distillation: Towards omni-supervised learning. In: *CVPR* (2018) [3](#), [4](#)
38. Rasmus, A., Berglund, M., Honkala, M., Valpola, H., Raiko, T.: Semi-supervised learning with ladder networks. In: *NeurIPS* (2015) [1](#), [5](#)
39. Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: Unified, real-time object detection. In: *CVPR* (2016) [4](#)
40. Redmon, J., Farhadi, A.: Yolo9000: better, faster, stronger. In: *CVPR* (2017) [4](#)
41. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: *NeurIPS* (2015) [2](#), [4](#), [6](#)
42. Rosenberg, C., Hebert, M., Schneiderman, H.: Semi-supervised self-training of object detection models. In: *Proc IEEE Workshops on Application of Computer Vision* (2005) [4](#)
43. RoyChowdhury, A., Chakrabarty, P., Singh, A., Jin, S., Jiang, H., Cao, L., Learned-Miller, E.: Automatic adaptation of object detectors to new domains using self-training. In: *CVPR* (2019) [2](#)
44. Sajjadi, M., Javanmardi, M., Tasdizen, T.: Mutual exclusivity loss for semi-supervised deep learning. In: *ICIP* (2016) [1](#), [5](#)
45. Sajjadi, M., Javanmardi, M., Tasdizen, T.: Regularization with stochastic transformations and perturbations for deep semi-supervised learning. In: *NeurIPS* (2016) [4](#)
46. Scudder, H.: Probability of error of some adaptive pattern-recognition machines. *IEEE Transactions on Information Theory* **11**(3) (1965) [2](#)
47. Shrivastava, A., Gupta, A., Girshick, R.: Training region-based object detectors with online hard example mining. In: *CVPR* (2016) [2](#), [15](#)
48. Simard, P.Y., Steinkraus, D., Platt, J.C., et al.: Best practices for convolutional neural networks applied to visual document analysis. In: *ICDAR* (2003) [1](#)
49. Sohn, K., Berthelot, D., Li, C.L., Zhang, Z., Carlini, N., Cubuk, E.D., Kurakin, A., Zhang, H., Raffel, C.: Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *arXiv preprint arXiv:2001.07685* (2020) [1](#), [2](#), [4](#), [5](#), [6](#), [7](#), [12](#), [22](#)
50. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: a simple way to prevent neural networks from overfitting. *JMLR* **15**(1), 1929–1958 (2014) [5](#)
51. Tan, M., Pang, R., Le, Q.V.: Efficientdet: Scalable and efficient object detection. *arXiv preprint arXiv:1911.09070* (2019) [4](#)
52. Tang, P., Ramaiah, C., Xu, R., Xiong, C.: Proposal learning for semi-supervised object detection. *arXiv preprint arXiv:2001.05086* (2020) [3](#), [4](#), [9](#)
53. Tang, Y., Wang, J., Gao, B., Dellandréa, E., Gaizauskas, R., Chen, L.: Large scale semi-supervised object detection using visual and semantic knowledge transfer. In: *CVPR* (2016) [2](#), [4](#)
54. Tarvainen, A., Valpola, H.: Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In: *NeurIPS* (2017) [1](#), [4](#), [5](#), [21](#)
55. Wang, K., Yan, X., Zhang, D., Zhang, L., Lin, L.: Towards human-machine cooperation: Self-supervised sample mining for object detection. In: *CVPR* (2018) [4](#)

56. Wu, Y., et al.: Tensorpack. <https://github.com/tensorpack/> (2016) 9
57. Xiao, J., Hays, J., Ehinger, K.A., Oliva, A., Torralba, A.: Sun database: Large-scale scene recognition from abbey to zoo. In: CVPR (2010) 12
58. Xie, Q., Dai, Z., Hovy, E., Luong, M.T., Le, Q.V.: Unsupervised data augmentation for consistency training. arXiv preprint arXiv:1904.12848 (2019) 1, 2, 4, 5, 7, 22
59. Xie, Q., Hovy, E., Luong, M.T., Le, Q.V.: Self-training with noisy student improves imagenet classification. arXiv preprint arXiv:1911.04252 (2019) 1, 2, 4, 5, 6, 22
60. Yarowsky, D.: Unsupervised word sense disambiguation rivaling supervised methods. In: 33rd annual meeting of the association for computational linguistics. pp. 189–196 (1995) 20
61. Yun, S., Han, D., Oh, S.J., Chun, S., Choe, J., Yoo, Y.: Cutmix: Regularization strategy to train strong classifiers with localizable features. In: ICCV (2019) 5
62. Zakharov, S., Kehl, W., Ilic, S.: Deceptionnet: Network-driven domain randomization. In: ICCV (2019) 2
63. Zhai, X., Oliver, A., Kolesnikov, A., Beyer, L.: S⁴l: Self-supervised semi-supervised learning. In: ICCV (2019) 4
64. Zhang, H., Cisse, M., Dauphin, Y.N., Lopez-Paz, D.: mixup: Beyond empirical risk minimization. arXiv preprint arXiv:1710.09412 (2017) 4, 5, 23
65. Zhang, Z., He, T., Zhang, H., Zhang, Z., Xie, J., Li, M.: Bag of freebies for training object detection neural networks. arXiv preprint arXiv:1902.04103 (2019) 5
66. Zhang, Z., Ringeval, F., Dong, B., Coutinho, E., Marchi, E., Schüller, B.: Enhanced semi-supervised learning for multimodal emotion recognition. In: ICASSP (2016) 15
67. Zhong, Z., Zheng, L., Kang, G., Li, S., Yang, Y.: Random erasing data augmentation. arXiv preprint arXiv:1708.04896 (2017) 4
68. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: ICCV (2017) 2
69. Zoph, B., Cubuk, E.D., Ghiasi, G., Lin, T.Y., Shlens, J., Le, Q.V.: Learning data augmentation strategies for object detection. arXiv preprint arXiv:1906.11172 (2019) 2, 3, 4, 5, 7, 8

A Learning Schedules

In this section, we provide complete descriptions on different learning schedules used in our experiments. Note that the schedule VOC is only used for experiments related to PASCAL VOC. Besides specified below, we adopt the learning settings as follows: <https://github.com/tensorpack/tensorpack/blob/master/examples/FasterRCNN/config.py>.

A.1 Quick

- **LR Decay:** $[0.01 (\leq 120k), 0.001 (\leq 160k), 0.0001 (\leq 180k)]$
- **Data processing:** Short edge size is sampled between 500 and 800 if the long edge is less than 1024 after resizing.
- **Batch per image for training Faster RCNN head:** 64

A.2 Standard, $[n] \times$

- **LR Decay:** $[0.01 (\leq 120k), 0.001 (\leq 160k), 0.0001 (\leq 180k)]$
- **LR Decay (2 \times):** $[0.01 (\leq 240k), 0.001 (\leq 320k), 0.0001 (\leq 360k)]$
- **LR Decay (3 \times):** $[0.01 (\leq 420k), 0.001 (\leq 500k), 0.0001 (\leq 540k)]$
- **Data processing:** Short edge size is fixed to 800 if the long edge is less than 1333 after resizing.
- **Batch per image for training Faster RCNN head:** 512

A.3 VOC

- **LR Decay:** $[0.001 (\leq 120k), 0.0005 (\leq 160k)]$
- **Data processing:** Short edge size is fixed to 600 if the long edge is less than 1000 after resizing. Image is resized to have its longer edge to be 1000 if long edge is longer than 1000.
- **Batch per image for training Faster RCNN head:** 256
- **RPN Anchor Sizes:** $[8, 16, 32]$

B Data Augmentation in STAC

This section provides more comprehensive results of Section 5.1 to validate the importance of data augmentation in STAC. In Table 6, we provide two rows of results with STAC (bottom) and the STAC without strong data augmentation, i.e., “Self-Training”. We observe significant gain in mAP on all cases, which validates the importance of the data augmentation in STAC.

Methods	5% COCO	10% COCO	100%
Self-Training	21.80±0.12	26.71±0.27	37.57
STAC	24.38±0.12	28.64±0.21	39.21

Table 6: Comparison in mAPs for different SSL methods on MS-COCO. We report the mean and standard deviation over 5 data folds for 5% and 10% protocols. “Self-Training” refers to STAC but without strong data augmentation on unlabeled data. We train STAC with the strong augmentation for unlabeled data.

C Extended Background: Unsupervised Loss in SSL

In this section, we extend Section 3.1 and provide unsupervised loss formulations for comprehensive list of SSL algorithms whose loss can be represented in Equation (1). For presentation clarity, let us reiterate definitions as follows:

$$\ell_u = \sum_{x \in \mathcal{X}} w(x) \ell(q(x), p(x)) \quad (6)$$

Here, we use $p(x)$ instead of $p(x; \theta)$ as in Equation (1) for generality. Instead, let us denote $p(x; \theta)$ as a prediction of the model with parameters θ at training.

Note that the unsupervised loss formulation of STAC is following the form of Noisy Student (Section C.9), which can be viewed as a combination of Self-Training (Section C.1) and strong data augmentation. While we have shown such a simple formulation of STAC brings in a significant performance gain at object detection, more complicated formulations (e.g., Mean Teacher (Section C.5) or MixMatch/ReMixMatch (Section C.10)) are amenable to be used in place of several design choices made for STAC. Further investigation of STAC variants is in the scope of the future work.

C.1 Bootstrapping (a.k.a. Self-Training) [60,32]

$$\ell(q, p) = \mathcal{H}(q, p) \quad (7)$$

$$w(x) = 1 \text{ if } \max(p(x; \tilde{\theta})) \geq \tau \text{ else } 0 \quad (8)$$

$$q(x) = p(x; \tilde{\theta}) \quad (9)$$

$$p(x) = p(x; \theta) \quad (10)$$

where $\tilde{\theta}$ is the parameter of the existing model, which usually refers to a model trained on labeled data only until convergence.

C.2 Entropy Minimization [16]

$$\ell(q, p) = \mathcal{H}(q, p) \quad (11)$$

$$w(x) = 1 \quad (12)$$

$$q(x) = p(x; \theta) \quad (13)$$

$$p(x) = p(x; \theta) \quad (14)$$

Note that gradient flows both to q and p . To our best knowledge, Entropy Minimization is the only method that backpropagates the gradient through q .

C.3 Pseudo Labeling [27]

$$\ell(q, p) = \mathcal{H}(q, p) \quad (15)$$

$$w(x) = 1 \text{ if } \max(p(x; \theta)) \geq \tau \text{ else } 0 \quad (16)$$

$$q(x) = \text{ONE_HOT}(\arg \max(p(x; \theta))) \quad (17)$$

$$p(x) = p(x; \theta) \quad (18)$$

C.4 Temporal Ensembling [25]

$$\ell(q, p) = \|q - p\|_2^2 \quad (19)$$

$$w(x) = 1 \quad (20)$$

$$q^{(t)}(x) = \alpha q^{(t-1)}(x) + (1 - \alpha)p(x; \theta) \quad (21)$$

$$p(x) = p(x; \theta) \quad (22)$$

We omit the ramp up and ramp down for $w(\cdot)$ in our formulation since it is dependent on the optimization framework. See [25] for more details.

C.5 Mean Teacher [54]

$$\ell(q, p) = \|q - p\|_2^2 \quad (23)$$

$$w(x) = 1 \quad (24)$$

$$q(x) = p(x; \theta^{\text{EMA}}), \theta^{\text{EMA}} = \alpha \theta^{\text{EMA}} + (1 - \alpha)\theta^{(t)} \quad (25)$$

$$p(x) = p(x; \theta^{(t)}) \quad (26)$$

We omit the ramp up and ramp down for $w(\cdot)$ in our formulation since it is dependent on the optimization framework. See [54] for more details.

C.6 Virtual Adversarial Training [35]

$$\ell(q, p) = \mathcal{H}(q, p) \quad (27)$$

$$w(x) = 1 \quad (28)$$

$$q(x) = p(x; \theta) \quad (29)$$

$$p(x) = p(\text{AP}(x); \theta), \text{AP}(\cdot): \text{adversarial perturbation} \quad (30)$$

C.7 Unsupervised Data Augmentation (UDA) [58]

UDA uses a weak ($\alpha(\cdot)$), such as translation and horizontal flip, to generate a pseudo label, and strong augmentation ($\mathcal{A}(\cdot)$), such as RandAugment [7] followed by Cutout [10], for model training.

$$\ell(q, p) = \mathcal{H}(q, p) \quad (31)$$

$$w(x) = 1 \text{ if } \max(p(\alpha(x); \theta)) \geq \tau \text{ else } 0 \quad (32)$$

$$q(x) \propto p(\alpha(x); \theta)^{\frac{1}{\tau}} \quad (33)$$

$$p(x) = p(\mathcal{A}(x); \theta) \quad (34)$$

C.8 FixMatch [49]

FixMatch also uses a weak ($\alpha(\cdot)$), such as translation and horizontal flip, to generate a pseudo label, and strong augmentation ($\mathcal{A}(\cdot)$), such as RandAugment [7] or CTAugment [3] followed by Cutout [10], for model training.

$$\ell(q, p) = \mathcal{H}(q, p) \quad (35)$$

$$w(x) = 1 \text{ if } \max(p(\alpha(x); \theta)) \geq \tau \text{ else } 0 \quad (36)$$

$$q(x) = \text{ONE_HOT}(\arg \max(p(\alpha(x); \theta))) \quad (37)$$

$$p(x) = p(\mathcal{A}(x); \theta) \quad (38)$$

C.9 Noisy Student [59]

$$\ell(q, p) = \mathcal{H}(q, p) \quad (39)$$

$$w(x) = 1 \text{ if } \max(p(x; \tilde{\theta})) \geq \tau \text{ else } 0 \quad (40)$$

$$q(x) = \text{ONE_HOT}(\arg \max(p(x; \tilde{\theta}))) \quad (41)$$

$$p(x) = p(\mathcal{A}(x); \theta) \quad (42)$$

where $\tilde{\theta}$ is the parameter of the model that is trained on labeled data only until convergence. In addition, Noisy Student perform data balancing across classes, which is not reflected in this formulation.

C.10 MixMatch [4]

Note that MixMatch uses MixUp [64] for unsupervised loss. It uses weak augmentation $\alpha(\cdot)$, such as translation and horizontal flip.

$$\ell(q, p) = \mathcal{H}(q, p) \quad (43)$$

$$w(x) = 1 \quad (44)$$

$$\tilde{q}(x) \propto \mathbb{E}_\alpha [p(\alpha(x); \theta)]^{\frac{1}{T}} \quad (45)$$

$$\tilde{q}(z) \propto \mathbb{E}_\alpha [p(\alpha(z); \theta)]^{\frac{1}{T}} \quad (46)$$

$$q(x) = \beta \tilde{q}(x) + (1 - \beta) \tilde{q}(z) \quad (47)$$

$$p(x) = p(\beta \alpha(x) + (1 - \beta) \alpha(z); \theta) \quad (48)$$

where x and z are unlabeled data and β is drawn from Beta distribution. While we present MixUp only between unlabeled data for presentation clarity, one may apply MixUp between labeled (with ground-truth label for \tilde{q}) and unlabeled data as well [4].

C.11 ReMixMatch [3]

Note that ReMixMatch uses MixUp [64] for unsupervised loss. It also uses weak augmentation $\alpha(\cdot)$, such as translation and horizontal flip, and strong augmentation $\mathcal{A}(\cdot)$, such as CTAugment [3].

$$\ell(q, p) = \mathbb{E}_{\mathcal{A}} [\mathcal{H}(q, p)] \quad (49)$$

$$w(x) = 1 \quad (50)$$

$$\tilde{q}(x) \propto p(\alpha(x); \theta)^{\frac{1}{T}} \quad (51)$$

$$\tilde{q}(z) \propto p(\alpha(z); \theta)^{\frac{1}{T}} \quad (52)$$

$$q(x) = \beta \tilde{q}(x) + (1 - \beta) \tilde{q}(z) \quad (53)$$

$$p(x) = p(\beta \mathcal{A}(x) + (1 - \beta) \mathcal{A}(z); \theta) \quad (54)$$

where x and z are unlabeled data and β is drawn from Beta distribution. While we present MixUp only between unlabeled data for presentation clarity, one may apply MixUp between labeled (with ground-truth label for \tilde{q}) and unlabeled data as well [3].